

Multimodal Data Mining in a Multimedia Database Based on Structured Max Margin Learning

ZHEN GUO and ZHONGFEI (MARK) ZHANG, SUNY Binghamton
ERIC P. XING and CHRISTOS FALOUTSOS, Carnegie Mellon University

Mining knowledge from a multimedia database has received increasing attentions recently since huge repositories are made available by the development of the Internet. In this article, we exploit the relations among different modalities in a multimedia database and present a framework for general multimodal data mining problem where image annotation and image retrieval are considered as the special cases. Specifically, the multimodal data mining problem can be formulated as a structured prediction problem where we learn the mapping from an input to the structured and interdependent output variables. In addition, in order to reduce the demanding computation, we propose a new max margin structure learning approach called Enhanced Max Margin Learning (EMML) framework, which is much more efficient with a much faster convergence rate than the existing max margin learning methods, as verified through empirical evaluations. Furthermore, we apply EMML framework to develop an effective and efficient solution to the multimodal data mining problem that is highly scalable in the sense that the query response time is independent of the database scale. The EMML framework allows an efficient multimodal data mining query in a very large scale multimedia database, and excels many existing multimodal data mining methods in the literature that do not scale up at all. The performance comparison with a state-of-the-art multimodal data mining method is reported for the real-world image databases.

Categories and Subject Descriptors: H.2.8 [**Database Management**]: Database Applications—*Data mining, Image databases*; H.3.3 [**Information Storage and Retrieval**]: Information Search and Retrieval—*Retrieval models*; I.5.1 [**Pattern Recognition**]: Models—*Structural*; J.3 [**Computer Applications**]: Life and Medical Sciences—*Biology and genetics*

General Terms: Algorithms, Experimentation

Additional Key Words and Phrases: Multimodal data mining, image annotation, image retrieval, max margin

ACM Reference Format:

Zhen Guo, Zhongfei (Mark) Zhang, Eric P. Xing, and Christos Faloutsos. 2016. Multimodal data mining in a multimedia database based on structured max margin learning. *ACM Trans. Knowl. Discov. Data* 10, 3, Article 23 (February 2016), 30 pages.

DOI: <http://dx.doi.org/10.1145/2742549>

1. INTRODUCTION

Mining knowledge from a multimedia database has received increasing attentions recently since huge repositories are made available by the development of the Internet. Multimedia data may consist of data in different modalities, such as digital images,

This work is supported in part by US NSF (IIS-0812114, CCF-1017828), the National Basic Research Program of China (2012CB316400), and Zhejiang Provincial Engineering Center on Media Data Cloud Processing and Analysis.

Authors' addresses: Z. Guo and Z. (Mark) Zhang, Computer Science Department, SUNY Binghamton, Binghamton, NY 13902; emails: {zguo, zhongfei}@cs.binghamton.edu; E. P. Xing and C. Faloutsos, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213; emails: {epxing, cfaloutsos}@cs.cmu.edu.

Permission to make digital or hard copies of part or all 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 show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2016 ACM 1556-4681/2016/02-ART23 \$15.00

DOI: <http://dx.doi.org/10.1145/2742549>

audio, video, and text data. In this context, a multimedia database refers to a data collection in which there are multiple modalities of data such as text and imagery. In this database system, the data in different modalities are related to each other. For example, text data are related to images as their annotation data. In this article, we focus on a multimedia database as an image database in which each image has a few textual words given as annotation.

Among different scenarios of the data mining from an image database, image retrieval plays an active role and the early research on the image retrieval problem is based on the content of images. The content-based approaches, however, suffer from a notorious bottleneck, semantic gap [Smeulders et al. 2000]. Recently, it is reported that this bottleneck may be reduced by the multimodal data mining approaches [Barnard et al. 2003; Feng et al. 2004], which take advantage of the fact that in many applications image data typically co-exist with other modalities of information such as text. The synergy between different modalities may be exploited to capture the high-level conceptual relationships. In addition, the multimodal data mining approaches bring more capabilities of knowledge discovery from multimedia database, such as the multimodal query consisting of multiple modalities of data.

Following the line of the multimodal data mining approach, in this article, we address a more general multimodal data mining problem in an image database as the problem of retrieving similar data and/or inferring new patterns to a multimodal query from the database. Specifically, in the context of this article, multimodal data mining refers to two aspects of activities. The first is the multimodal retrieval. This is the scenario where a multimodal query consisting of either textual words alone, or imagery alone, or in any combination is entered and an expected retrieved data modality is specified that can also be text alone, or imagery alone, or in any combination; the retrieved data based on a pre-defined similarity criterion are returned back to the user. The image annotation and image retrieval are considered as the special cases in this scenario.

The second is the multimodal inference. While the retrieval-based multimodal data mining has its standard definition in terms of the semantic similarity between the query and the retrieved data from the database, the inference-based mining depends on the specific applications. In this article, we investigate the application of the fruit fly image database mining. Consequently, the inference-based multimodal data mining may include many different scenarios. A typical scenario is the multimodal inference between different stages. There are many interesting questions a biologist may want to ask in the fruit fly research given such a multimodal mining capability. For example, given an embryo image in stage 5, what is the corresponding image in stage 7 for an image-to-image three-stage inference? What is the corresponding annotation for this image in stage 7 for an image-to-word three-stage inference? The multimodal mining technique we have developed in this article also provides this type of inference capability, in addition to the multimodal retrieval capability.

In order to exploit the synergy within the multimodal data, the relationships among the different modalities need to be learned. For an image database, we need to learn the relationship between images and text. The learned relationship between images and text can then be further used in multimodal data mining. Without loss of generality, we start with a special case of the multimodal data mining problem—image annotation, where the input is an image query and the expected output is the annotation words. We show later that this approach is also valid to the general multimodal data mining problem defined in this article. The image annotation problem can be formulated as a structured prediction problem where the input (image) \mathbf{x} and the output (annotation) \mathbf{y} are structures and the goal is to obtain the mapping between these two structured spaces. Specifically, an image can be partitioned into a set of blocks which form a structure in the sense that the adjacent blocks are more close to each other than the

non-adjacent blocks. On the other hand, depending on the context of the image, some certain words are more likely to occur together as the annotation than other words, thus, the words are interdependent on each other and form a structured space which can be represented by a vector space with each entry for one word. For example, in the context of landscape, “sea” and “ship” are likely to appear together as the annotation whereas “sea” and “car” are very unlikely to occur together. Under this setting, the learning task is, therefore, formulated as finding a function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ such that

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}) \quad (1)$$

is the desired output for any input \mathbf{x} .

It is worth mentioning that it is intractable to solve the structured prediction problem by the multi-class approach because the number of the structures is exponential with respect to the size of the structured space. Another challenge in the structured prediction problem is that the number of the constraints is increased combinatorially in terms of the size of the structured space. Tsochantaridis et al. [2004] propose a cutting plane algorithm to reduce the number of the constraints. However, it needs to compute the most violated constraint, which involve another optimization problem in the output space, and thus, requires additional computation. In this article, we propose a max margin learning approach on the structured output space called Enhanced Max Margin Learning (EMML), which substantially reduces the number of constraints by taking advantage of the redundancy among the constraints, and thus, provides a much more efficiency with a much faster convergence rate than the existing approaches.

Note that the proposed approach is general that can be applied to any structured prediction problems. In this article, EMML is applied to the multimodal data mining from image databases with textual annotations as the special cases of multimedia databases. In order to truly capture the difficulties in real scenarios such as Web image retrieval and biological image analysis, we evaluate the proposed framework on two different real-world image databases including an image database from various crawled Web pages and the Berkeley Drosophila embryo image database. Extensive empirical evaluations against a state-of-the-art method on these databases are reported.

In summary, we highlight the contributions of this work as follows.

- We address the general multimodal data mining problem in a multimedia database, which includes the traditional multimodal retrieval and the multimodal inference.
- We propose an efficient approach called EMML to solve the structured learning problem. In EMML, the number of the constraints is dramatically reduced by exploiting the redundancy among the constraints, and thus, the proposed algorithm converges much faster than the existing approaches.
- we have applied the EMML approach to develop an effective and efficient solution to the multimodal data mining problem that is highly scalable in the sense that the query response time is independent of the database scale. The EMML framework allows an efficient multimodal data mining query in a very large scale multimedia database, and excels many existing multimodal data mining methods in the literature that do not scale up at all. This advantage is also supported through the complexity analysis as well as empirical evaluations against a state-of-the-art multimodal data mining method from the literature.

The rest of the article is organized as follows. Section 2 discusses the related work on multimodal approaches and the max margin learning in the structured output space. In Section 3, the learning in the structured output space is described using the example of the image annotation and the general EMML framework is proposed. The solution to the problem of multimodal data mining is given based on the EMML framework

in Section 4. The extensive experiments to evaluate the proposed EML framework against a state-of-the-art multimodal data mining method are reported in Section 5. We conclude the article in Section 6. This article extends [Guo et al. 2007a] with additional theoretical analysis and complexity analysis as well as empirical results.

2. RELATED WORK

Multimodal approaches have recently received substantial attentions since Barnard et al. and Duygulu et al. started their pioneering work on image annotation [Duygulu et al. 2002; Barnard et al. 2003]. Recently, there have been many studies [Blei and Jordan 2003; Chang et al. 2003; Feng et al. 2004; Pan et al. 2004; Wu et al. 2005; Datta et al. 2006] on the multimodal approaches. Blei and Jordan [2003] present a correspondence latent Dirichlet allocation (CORR-LDA) model, which implicitly establishes the correspondence between the image regions and annotation words by assuming that the regions of the image can be conditional on any ensemble of factors but the annotation words must be conditional on factors that are present in the image. CORR-LDA models the annotated data from the probability point of view while EML is an approach based on the max margin learning.

The learning with structured output variables covers many natural learning tasks including named entity recognition, natural language parsing, and label sequence learning. Recently the problem of image annotation and image retrieval has been formulated as a structured learning problem [Guo et al. 2007b]. The challenge of learning with structured output variables is that the number of the structures is exponential in terms of the size of the structured output space. Thus, the problem is intractable if we treat each structure as a separate class. Consequently, the conventional multiclass learning approaches are not well fitted into the learning with structured output variables. There have been many studies on the structured model that include conditional random fields [Lafferty et al. 2001], maximum entropy model [McCallum et al. 2000], graph model [Chu et al. 2004], semi-supervised learning [Brefeld and Scheffer 2006], and max margin approaches [Altun et al. 2003; Taskar et al. 2003; Tsochantaridis et al. 2004; Daume III and Marcu 2005]. Lafferty et al. [2001] propose the Conditional Random Fields (CRF) model for labeling sequence data. CRF is a probabilistic model that assumes that the data sequence and the label sequence have the same dimension, for example, in the part-of-speech tagging problem where each component (tag) in the label sequence corresponds to a word in the sentence. The image annotation problem as well as the general multimodal data mining problem considered in this article, however, obviously violates this assumption since an image might have several annotation words. Therefore, the CRF model is not applicable to the problems addressed in this article.

As an effective approach to learning with structured output variables, the max margin principle has received substantial attentions since it is formulated in the framework of the support vector machine (SVM) [Vapnik 1995]. In addition, the perceptron algorithm is also used to explore the max margin classification [Freund and Schapire 1999]. Taskar et al. [2005] reduce the number of the constraints by considering the dual of the loss-augmented problem. However, the number of the constraints in their approach is still intractably large for a large structured output space and a large training set.

For learning with structured output variables, Tsochantaridis et al. [2004] propose a cutting plane algorithm which finds a small set of active constraints. One issue of this algorithm is that it needs to compute the most violated constraint which would involve another optimization problem in the output space. In EML, instead of selecting the most violated constraint, we arbitrarily select a constraint that violates the optimality condition of the optimization problem. Thus, the selection of the constraints does not

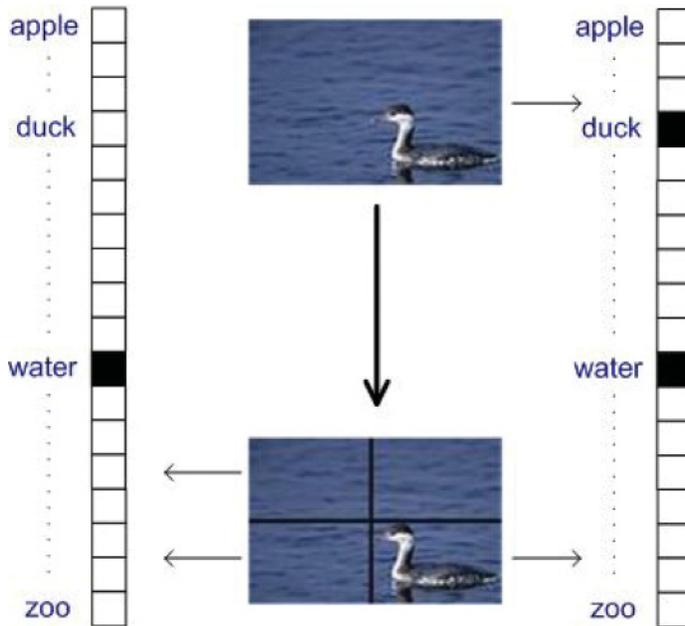


Fig. 1. An illustration of the image partitioning and the structured output word space.

involve any optimization problem. Osuna et al. [1997] propose the decomposition algorithm for the SVM. There are many efficient algorithms on the SVMs in the literature. Osuna et al.'s decomposition algorithm, however, is unique in the sense that it makes use of the fact that only the support vectors determine the classification function and the removal of the other examples does not affect the solution. In other words, Osuna et al.'s approach incorporates the essential features of the SVMs. Therefore, Osuna et al.'s approach achieves the promising performance. In EMML, we extend their idea to the scenario of learning with structured output variables.

3. EMML FRAMEWORK

We use image annotation problem to derive the EMML framework; the framework, however, is applicable to general multimodal data mining problems. Assume that the image database consists of a set of instances $S = \{(I_i, W_i)\}_{i=1}^L$, where each instance consists of an image object I_i and the corresponding annotation word set W_i . First, we partition an image into a set of blocks. Thus, an image can be represented by a set of sub-images. The feature vector in the feature space for each block can be computed from the selected feature representation. Consequently, an image is represented as a set of feature vectors in the feature space. A clustering algorithm is then applied to the whole feature space to group similar feature vectors together. The centroid of a cluster represents a visual representative (we refer it to VRep in this article) in the image space. An illustrative example is shown in Figure 1 where there are two VReps, *water* and *duck* in the water. The corresponding annotation word set can be easily obtained for each VRep. Consequently, the image database becomes the VRep-word pairs $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, where n is the number of the clusters, \mathbf{x}_i is a VRep object, and \mathbf{y}_i is the word annotation set corresponding to this VRep object. Another simple method to obtain the VRep-word pairs is that we randomly select some images from the image database and each image is viewed as a VRep.

3.1. Learning in the Structured Output Space

Suppose that there are W distinct annotation words. An arbitrary subset of annotation words is represented by a binary vector $\bar{\mathbf{y}}$ whose length is W ; the j th component $\bar{y}_j = 1$ if the j th word occurs in this subset, and 0 otherwise. All possible binary vectors form the word space \mathcal{Y} . We use \mathbf{w}_j to denote the vector representation of the j th word in the whole word set. We use \mathbf{x} to denote an arbitrary vector in the feature space. In the illustrative example Figure 1, the original image is annotated by *duck* and *water*, which are represented by a binary vector. There are two VReps after the clustering and each has a different annotation. In the word space, a word may be related to other words. For example, *duck* and *water* are related to each other because *water* is more likely to occur when *duck* is one of the annotation words. Consequently, the annotation word space is a structured output space where the elements are interdependent on each other.

The relationship between the input example VRep \mathbf{x} and an arbitrary output $\bar{\mathbf{y}}$ is represented as the joint feature mapping $\Phi(\mathbf{x}, \bar{\mathbf{y}})$, $\Phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$, where d is the dimension of the joint feature space. It can be expressed as a linear combination of the joint feature mapping between \mathbf{x} and all the unit vectors. That is

$$\Phi(\mathbf{x}, \bar{\mathbf{y}}) = \sum_{j=1}^W \bar{y}_j \Phi(\mathbf{x}, \mathbf{e}_j) \quad (2)$$

where \mathbf{e}_j is the j th unit vector. The score between \mathbf{x} and $\bar{\mathbf{y}}$ can be expressed as a linear combination of each component in the joint feature representation: $f(\mathbf{x}, \bar{\mathbf{y}}) = \langle \alpha, \Phi(\mathbf{x}, \bar{\mathbf{y}}) \rangle$. Then, the learning task is to find the optimal weight vector α such that the prediction error is minimized for all the training instances. That is

$$\arg \max_{\bar{\mathbf{y}} \in \mathcal{Y}_i} f(\mathbf{x}_i, \bar{\mathbf{y}}) \approx \mathbf{y}_i, \quad i = 1, \dots, n$$

where

$$\mathcal{Y}_i = \text{perm}(\mathbf{y}_i) \quad (3)$$

and $\text{perm}(\mathbf{y}_i)$ denotes the set that consists of all the vectors transformed from \mathbf{y}_i by permutation. The space \mathcal{Y}_i includes all the vectors that have the same number of 1 as \mathbf{y}_i . So, the prediction includes the same number of words as the label \mathbf{y}_i , which is a reasonable assumption since the label \mathbf{y}_i is the best prediction. The advantage of this assumption is that the size of space \mathcal{Y}_i is not very large and the computation of loss function $l(\bar{\mathbf{y}}, \mathbf{y}_i)$ is very convenient. We use $\Phi_i(\bar{\mathbf{y}})$ to denote $\Phi(\mathbf{x}_i, \bar{\mathbf{y}})$. To make the prediction to be the true output \mathbf{y}_i , it is obvious to have

$$\alpha^\top \Phi_i(\mathbf{y}_i) \geq \alpha^\top \Phi_i(\bar{\mathbf{y}}) \quad \forall \bar{\mathbf{y}} \in \mathcal{Y}_i \setminus \{\mathbf{y}_i\}$$

where $\mathcal{Y}_i \setminus \{\mathbf{y}_i\}$ denotes the removal of the element \mathbf{y}_i from the set \mathcal{Y}_i . In order to accommodate the prediction error on the training examples, we introduce the slack variable ξ_i for each training sample. The above constraint then becomes

$$\alpha^\top \Phi_i(\mathbf{y}_i) \geq \alpha^\top \Phi_i(\bar{\mathbf{y}}) - \xi_i, \quad \xi_i \geq 0 \quad \forall \bar{\mathbf{y}} \in \mathcal{Y}_i \setminus \{\mathbf{y}_i\}.$$

The prediction error on the training instances is measured by the loss function, which is the distance between the true output \mathbf{y}_i and the prediction $\bar{\mathbf{y}}$. The loss function measures the goodness of the learning model. The standard zero-one classification loss is not suitable for the structured output space since it cannot capture the actual difference from the true output. Consequently, it is natural to define the loss function $l(\bar{\mathbf{y}}, \mathbf{y}_i)$ as the number of the different entries in these two vectors. We then include the

loss function in the constraints as is proposed by Taskar et al. [2005]

$$\boldsymbol{\alpha}^\top \Phi_i(\mathbf{y}_i) \geq \boldsymbol{\alpha}^\top \Phi_i(\bar{\mathbf{y}}) + l(\bar{\mathbf{y}}, \mathbf{y}_i) - \xi_i. \quad (4)$$

We interpret $\frac{1}{\|\boldsymbol{\alpha}\|} \boldsymbol{\alpha}^\top [\Phi_i(\mathbf{y}_i) - \Phi_i(\bar{\mathbf{y}})]$ as the margin of \mathbf{y}_i over another $\bar{\mathbf{y}} \in \mathcal{Y}^{(i)}$. We then rewrite the above constraint as $\frac{1}{\|\boldsymbol{\alpha}\|} \boldsymbol{\alpha}^\top [\Phi_i(\mathbf{y}_i) - \Phi_i(\bar{\mathbf{y}})] \geq \frac{1}{\|\boldsymbol{\alpha}\|} [l(\bar{\mathbf{y}}, \mathbf{y}_i) - \xi_i]$. Thus, minimizing $\|\boldsymbol{\alpha}\|$ maximizes such margin.

The goal now is to solve the optimization problem

$$\begin{aligned} \min \quad & \frac{1}{2} \|\boldsymbol{\alpha}\|^2 + C \sum_{i=1}^n \xi_i^r \\ \text{s.t.} \quad & \boldsymbol{\alpha}^\top \Phi_i(\mathbf{y}_i) \geq \boldsymbol{\alpha}^\top \Phi_i(\bar{\mathbf{y}}) + l(\bar{\mathbf{y}}, \mathbf{y}_i) - \xi_i \\ & \forall \bar{\mathbf{y}} \in \mathcal{Y}_i \setminus \{\mathbf{y}_i\}, \quad \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned} \quad (5)$$

where $r = 1, 2$ corresponds to the linear or quadratic slack variable penalty. In this article, we use the linear slack variable penalty. For $r = 2$, we obtain similar results. $C > 0$ is a constant that controls the tradeoff between the training error minimization and the margin maximization.

The justification for Equation (4) is given by the following proposition.

PROPOSITION 1. *Denote by $(\boldsymbol{\alpha}^*, \boldsymbol{\xi}^*)$ the optimal solution to Equation (5), then $\frac{1}{n} \sum_{i=1}^n (\boldsymbol{\alpha}^{*\top} \Phi_i(\mathbf{y}_i) + \xi_i^*)$ is an upper bound on the empirical risk $\mathcal{R}(f)$.*

PROOF. According to the definition, the empirical risk $\mathcal{R}(f)$ is given by

$$\mathcal{R}(f) = \frac{1}{n} \sum_{i=1}^n l(\mathbf{y}_i, f(\mathbf{x}_i)) \leq \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\alpha}^{*\top} \Phi_i(\mathbf{y}_i) + \xi_i^*). \quad \square \quad (6)$$

Note that in the above formulation, we do not introduce the relationships between different words in the word space. However, the relationships between different words are implicitly included in the VRep-word pairs because the related words are more likely to occur together. Thus, Equation (5) is in fact a structured optimization problem.

The above discussion assumes that the word space \mathcal{Y} only contains the binary vector $\bar{\mathbf{y}}$. Our model, however, can be easily extended to handle the general vector $\bar{\mathbf{y}}$ with the continuous values. In this case, each entry in $\bar{\mathbf{y}}$ is a real number, which indicates the relevant degree associated with the image. This problem can be considered as the regression problem in the structured output space. Since Equations (2) and (3) are defined over the general vectors, the optimization problem Equation (5) can be used to solve the regression problem in the structured output space. Consequently, Equation (1) can be applied to compute the relevant degree between the images and the words.

3.2. Optimality Conditions

One can solve the optimization problem Equation (5) in the primal space—the space of the parameters $\boldsymbol{\alpha}$. In fact this problem is intractable when the structured output space is large because the number of the constraints is increased combinatorially in terms of the size of the output space. As in the traditional SVM, the solution can be obtained by solving this quadratic optimization problem in the dual space—the space of the Lagrange multipliers. Vapnik [1995] and Boyd and Vandenberghe [2004] have an excellent review for the related optimization problems.

The dual problem formulation has an important advantage over the primal problem: it only depends on the inner products in the joint feature representation defined by Φ , allowing the use of a kernel function. We introduce the Lagrange multiplier $\mu_{i,\bar{\mathbf{y}}}$ for each constraint to form the Lagrangian. We define $\Phi_{i,\mathbf{y}_i,\bar{\mathbf{y}}} = \Phi_i(\mathbf{y}_i) - \Phi_i(\bar{\mathbf{y}})$ and the kernel

function $K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_j, \bar{\mathbf{y}})) = \langle \Phi_{i, \mathbf{y}_i, \bar{\mathbf{y}}}, \Phi_{j, \mathbf{y}_j, \bar{\mathbf{y}}} \rangle$. The derivatives of the Lagrangian over α and ξ_i should be equal to zero. Substituting these conditions into the Lagrangian, we obtain the following Lagrange dual problem

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{\substack{i,j \\ \bar{\mathbf{y}} \neq \mathbf{y}_i \\ \bar{\mathbf{y}} \neq \mathbf{y}_j}} \mu_{i, \bar{\mathbf{y}}} \mu_{j, \bar{\mathbf{y}}} K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_j, \bar{\mathbf{y}})) - \sum_{\substack{i \\ \bar{\mathbf{y}} \neq \mathbf{y}_i}} \mu_{i, \bar{\mathbf{y}}} l(\bar{\mathbf{y}}, \mathbf{y}_i) \\ \text{s.t.} \quad & \sum_{\bar{\mathbf{y}} \neq \mathbf{y}_i} \mu_{i, \bar{\mathbf{y}}} \leq C \quad \mu_{i, \bar{\mathbf{y}}} \geq 0 \quad i = 1, \dots, n. \end{aligned} \quad (7)$$

After this dual problem is solved, the solution to the primal problem Equation (5) can be obtained by

$$\alpha = \sum_{i, \bar{\mathbf{y}}} \mu_{i, \bar{\mathbf{y}}} \Phi_{i, \mathbf{y}_i, \bar{\mathbf{y}}}. \quad (8)$$

For each training example, there are a number of constraints related to it. We use the subscript i to represent the part related to the i th example in the matrix. For example, let μ_i be the vector with entries $\mu_{i, \bar{\mathbf{y}}}$. We stack the μ_i together to form the vector μ . That is, $\mu = [\mu_1^\top \cdots \mu_n^\top]^\top$. Similarly, let \mathbf{S}_i be the vector with entries $l(\bar{\mathbf{y}}, \mathbf{y}_i)$. We stack \mathbf{S}_i together to form the vector \mathbf{S} . That is, $\mathbf{S} = [\mathbf{S}_1^\top \cdots \mathbf{S}_n^\top]^\top$. The lengths of μ and \mathbf{S} are the same. We define \mathbf{A}_i as the vector that has the same length as that of μ , where only the entries in the vector \mathbf{A}_i corresponding to μ_i are 1 and all other entries are 0. So, $\mathbf{A}_i^\top \mu = \sum_{\bar{\mathbf{y}} \neq \mathbf{y}_i} \mu_{i, \bar{\mathbf{y}}}$, which is exactly the left hand of the constraints in the optimization problem Equation (7). Let $\mathbf{A} = [\mathbf{A}_1 \cdots \mathbf{A}_n]^\top$. Let matrix \mathbf{D} represent the kernel matrix where each entry is $K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_j, \bar{\mathbf{y}}))$. Let \mathbf{C} be the vector where each entry is constant C .

With the above notations, we rewrite the Lagrange dual problem as follows:

$$\begin{aligned} \min \quad & \frac{1}{2} \mu^\top \mathbf{D} \mu - \mu^\top \mathbf{S} \\ \text{s.t.} \quad & \mathbf{A} \mu \leq \mathbf{C} \\ & \mu \geq 0 \end{aligned} \quad (9)$$

where \leq and \geq represent the vector comparison defined as entry-wise less than or equal to and greater than or equal to, respectively.

Equation (9) has the same number of the constraints as Equation (5). However, in Equation (9), most of the constraints are lower bound constraints ($\mu \geq 0$), which define the feasible region. Other than these lower bound constraints, the rest of the constraints determine the complexity of the optimization problem. Therefore, the number of the constraints is considered to be reduced in Equation (9). However, the challenge still exists to solve it efficiently since the number of the dual variables is still huge.

The optimality conditions of Equation (9) need to be considered to develop an efficient algorithm. The Lagrangian of Equation (9) can be obtained:

$$L = L(\mu, \gamma, \pi) = \frac{1}{2} \mu^\top \mathbf{D} \mu - \mu^\top \mathbf{S} - \gamma^\top (\mathbf{C} - \mathbf{A} \mu) - \pi^\top \mu$$

where γ and π are the non-negative Lagrange multipliers. Since the optimization problem Equation (9) is a convex optimization problem, we have the following

Karush–Kuhn–Tucker (KKT) conditions [Boyd and Vandenberghe 2004] of optimality.

$$\begin{aligned}
\frac{\partial L}{\partial \boldsymbol{\mu}} &= \mathbf{D}\boldsymbol{\mu} - \mathbf{S} + \mathbf{A}^\top \boldsymbol{\gamma} - \boldsymbol{\pi} = 0 \\
\boldsymbol{\gamma}^\top (\mathbf{C} - \mathbf{A}\boldsymbol{\mu}) &= 0 \\
\boldsymbol{\pi}^\top \boldsymbol{\mu} &= 0 \\
\boldsymbol{\gamma} &\geq 0 \\
\boldsymbol{\pi} &\geq 0 \\
\mathbf{A}\boldsymbol{\mu} &\leq \mathbf{C} \\
\boldsymbol{\mu} &\geq 0.
\end{aligned} \tag{10}$$

For the optimization problem Equation (9), the KKT conditions provide necessary and sufficient conditions for optimality.

Based on the KKT conditions, we consider the two possible values that $\sum_{\bar{y} \neq y_i} \mu_{i,\bar{y}}$ can have. Note that $\sum_{\bar{y} \neq y_i} \mu_{i,\bar{y}} = \mathbf{A}_i^\top \boldsymbol{\mu}$.

(1) $\mathbf{A}_i^\top \boldsymbol{\mu} < C$:

According to the second equation in Equation (10), we have $\gamma_i = 0$. Consequently, the first equation in Equation (10) suggests that $(\mathbf{D}\boldsymbol{\mu})_{i,\bar{y}} - \mathbf{S}_{i,\bar{y}} = \pi_{i,\bar{y}}$ holds true, leading to $(\mathbf{D}\boldsymbol{\mu})_{i,\bar{y}} - \mathbf{S}_{i,\bar{y}} \geq 0$ since $\pi_{i,\bar{y}} \geq 0$. In other words, the following two conditions do not hold true simultaneously if $\boldsymbol{\mu}$ is the optimal solution to the optimization problem Equation (9).

$$\begin{aligned}
\mathbf{A}_i^\top \boldsymbol{\mu} &< C \\
(\mathbf{D}\boldsymbol{\mu})_{i,\bar{y}} - \mathbf{S}_{i,\bar{y}} &< 0.
\end{aligned} \tag{11}$$

In the Section 3.3, Equation (13) [which is consistent with Equation (11)] is used to prove the correctness of the decomposition algorithm proposed in this article, as shown in Theorem 3.2.

(2) $\mathbf{A}_i^\top \boldsymbol{\mu} = C$:

In this case, the second equation in Equation (10) holds true no matter what values $\boldsymbol{\gamma}$ take. In fact, $\boldsymbol{\gamma}$ is determined by the first equation. In other words, there always exist the appropriate $\boldsymbol{\gamma}, \boldsymbol{\pi}$, which meet the KKT conditions.

Incorporating these optimal conditions into the optimization problem leads to an efficient algorithm. To achieve this goal, we need to derive an algorithm to use this information.

3.3. Decomposition Algorithm

Osuna et al. [1997] propose a decomposition algorithm for the SVM learning over large datasets. We extend this idea to learning with the structured output space. We decompose the constraints of the optimization problem Equation (5) into two sets: the working set B and the nonactive set N. The Lagrange multipliers are also correspondingly partitioned into two parts $\boldsymbol{\mu}_B$ and $\boldsymbol{\mu}_N$. We are interested in the subproblem defined only for the dual variable set $\boldsymbol{\mu}_B$ when keeping $\boldsymbol{\mu}_N = 0$.

This subproblem is formulated as follows:

$$\begin{aligned}
\min \quad & \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{D}\boldsymbol{\mu} - \boldsymbol{\mu}^\top \mathbf{S} \\
\text{s.t.} \quad & \mathbf{A}\boldsymbol{\mu} \leq \mathbf{C} \\
& \boldsymbol{\mu}_B \geq 0, \quad \boldsymbol{\mu}_N = 0.
\end{aligned} \tag{12}$$

It is clearly true that we can move those $\mu_{i,\bar{\mathbf{y}}} = 0, \mu_{i,\bar{\mathbf{y}}} \in \mu_B$ to set μ_N without changing the objective function. Furthermore, we can move those $\mu_{i,\bar{\mathbf{y}}} \in \mu_N$ satisfying certain conditions to set μ_B to form a new optimization subproblem, which yields a strict decrease in the objective function in Equation (9) when the new subproblem is optimized. This property is guaranteed by the following theorem.

THEOREM 3.1. *Given an optimal solution to the subproblem defined on μ_B in Equation (12), if the following conditions hold true:*

$$\begin{aligned} \exists i, \quad \sum_{\bar{\mathbf{y}}} \mu_{i,\bar{\mathbf{y}}} &< C \\ \exists \mu_{i,\bar{\mathbf{y}}} \in \mu_N, \quad \alpha^\top \Phi_{i,\mathbf{y}_i,\bar{\mathbf{y}}} - l(\bar{\mathbf{y}}, \mathbf{y}_i) &< 0 \end{aligned} \quad (13)$$

the operation of moving the Lagrange multiplier $\mu_{i,\bar{\mathbf{y}}}$ satisfying Equation (13) from set μ_N to set μ_B generates a new optimization subproblem that yields a strict decrease in the objective function in Equation (9) when the new subproblem in Equation (12) is optimized.

PROOF. Suppose that the current optimal solution is μ . Let δ be a small positive number. Let $\bar{\mu} = \mu + \delta e_r$, where e_r is the r th unit vector and $r = (i, \bar{\mathbf{y}})$ denotes the Lagrange multiplier satisfying condition Equation (13). Thus, the objective function becomes

$$\begin{aligned} \mathbf{W}(\bar{\mu}) &= \frac{1}{2}(\mu + \delta e_r)^\top \mathbf{D}(\mu + \delta e_r) - (\mu + \delta e_r)^\top \mathbf{S} \\ &= \frac{1}{2}(\mu^\top \mathbf{D}\mu + \delta e_r^\top \mathbf{D}\mu + \delta \mu^\top \mathbf{D}e_r + \delta^2 e_r^\top \mathbf{D}e_r) \\ &\quad - \mu^\top \mathbf{S} - \delta e_r^\top \mathbf{S} \\ &= \mathbf{W}(\mu) + \frac{1}{2}(\delta e_r^\top \mathbf{D}\mu + \delta \mu^\top \mathbf{D}e_r + \delta^2 e_r^\top \mathbf{D}e_r) \\ &\quad - \delta e_r^\top \mathbf{S} \\ &= \mathbf{W}(\mu) + \delta e_r^\top \mathbf{D}\mu - \delta e_r^\top \mathbf{S} + \frac{1}{2}\delta^2 e_r^\top \mathbf{D}e_r \\ &= \mathbf{W}(\mu) + \delta(\alpha^\top \Phi_{i,\mathbf{y}_i,\bar{\mathbf{y}}} - l(\bar{\mathbf{y}}, \mathbf{y}_i)) + \frac{1}{2}\delta^2 \|\Phi_{i,\mathbf{y}_i,\bar{\mathbf{y}}}\|^2. \end{aligned}$$

Since $\alpha^\top \Phi_{i,\mathbf{y}_i,\bar{\mathbf{y}}} - l(\bar{\mathbf{y}}, \mathbf{y}_i) < 0$, for small enough δ , we have $\mathbf{W}(\bar{\mu}) < \mathbf{W}(\mu)$. For small enough δ , the constraints $\mathbf{A}\bar{\mu} \leq \mathbf{C}$ are also valid. Therefore, when the new optimization subproblem in Equation (12) is optimized, there must be an optimal solution no worse than $\bar{\mu}$. \square

Note that Equations (11) and (13) are consistent with each other although they are derived in different scenarios. From Equation (8), we have

$$\begin{aligned} \alpha^\top \Phi_{i,\mathbf{y}_i,\bar{\mathbf{y}}} - l(\bar{\mathbf{y}}, \mathbf{y}_i) &= \sum_{j,\bar{\mathbf{y}}} \mu_{j,\bar{\mathbf{y}}} \langle \Phi_{j,\mathbf{y}_j,\bar{\mathbf{y}}}, \Phi_{i,\mathbf{y}_i,\bar{\mathbf{y}}} \rangle - l(\bar{\mathbf{y}}, \mathbf{y}_i) \\ &= (\mathbf{D}\mu)_{i,\bar{\mathbf{y}}} - \mathbf{S}_{i,\bar{\mathbf{y}}}. \end{aligned}$$

In fact, the optimal solution is obtained when there is no Lagrange multiplier satisfying the condition Equation (13). This is guaranteed by the following theorem.

THEOREM 3.2. *The optimal solution to the optimization problem in Equation (9) is achieved if and only if the condition Equation (13) does not hold true.*

PROOF. If the optimal solution $\hat{\mu}$ is achieved, the condition Equation (13) must not hold true. Otherwise, $\hat{\mu}$ is not optimal according to Theorem 3.1. To prove in the reverse direction, we consider the KKT conditions Equation (10).

For the optimization problem Equation (9), the KKT conditions provide necessary and sufficient conditions for optimality. One can check that the condition Equation (13) violates the KKT conditions. On the other hand, one can check that the KKT conditions are satisfied when condition Equation (13) does not hold true according to the optimality condition analysis in Section 3.2. Therefore, the optimal solution is achieved when condition Equation (13) does not hold true. \square

According to the above theorems, we propose the EMLL algorithm listed in Algorithm 1.

ALGORITHM 1: EMLL Algorithm

Input: n labeled examples, dual variable set μ .

Output: Optimized μ

1: **procedure**

2: Arbitrarily decompose μ into two sets: μ_B and μ_N .

3: Solve the subproblem in Equation (12) defined by the variables in μ_B .

4: While there exists $\mu_{i,\bar{y}} \in \mu_B$ such that $\mu_{i,\bar{y}} = 0$, move it to set μ_N .

5: While there exists $\mu_{i,\bar{y}} \in \mu_N$ satisfying condition Equation (13), move it to set μ_B .
If no such $\mu_{i,\bar{y}} \in \mu_N$ exists, the iteration exits.

6: Goto Step 3.

7: **end procedure**

Note that in Step 5, we only need to find one dual variable satisfying Equation (13). We need to examine all the dual variables in the set μ_N only when no dual variable satisfies Equation (13). It is fast to examine the dual variables in the set μ_N even if the number of the dual variables is large.

3.4. Analysis

In addition to the conditions of selecting the constraints as shown in the above theorems, we are also interested in how fast the algorithm converges. The convergence rate and computation complexity of the algorithm can be derived based on the following Lemma.

LEMMA 3.3. *Let \mathbf{J} be a nonnegative symmetric matrix with the positive diagonal entries and assume that we have the following quadratic problem:*

$$P(\mathbf{v}) = \frac{1}{2} \mathbf{v}^\top \mathbf{J} \mathbf{v} - \mathbf{h}^\top \mathbf{v}$$

where $\mathbf{v} \geq 0$ with $v_r = 0$ and \mathbf{h} is a vector with appropriate size. Assume that $h_r \geq \sum_s v_s J_{rs}$. Minimizing $P(\mathbf{v})$ with respect to v_r while keeping all other entries fixed will decrease the objective by

$$\frac{(h_r - \sum_s v_s J_{rs})^2}{2J_{rr}}.$$

PROOF. The terms containing v_r are

$$P(v_r) = \frac{v_r^2}{2} J_{rr} - v_r \left(h_r - \sum_s v_s J_{rs} \right).$$

Setting the derivative over v_r to zero leads to

$$v_r^* = \frac{h_r - \sum_s v_s J_{rs}}{J_{rr}}$$

$$P(v_r^*) = -\frac{(h_r - \sum_s v_s J_{rs})^2}{2J_{rr}}.$$

Note that $v_r^* \geq 0$ since $h_r \geq \sum_s v_s J_{rs}$ and $J_{rr} > 0$. \square

The lower bound of decrease in the objective function achieved by Step 3 in Algorithm 1 can be obtained by Lemma 3.3, as shown in the following theorem.

THEOREM 3.4. *Let $R_i = \max_{\bar{\mathbf{y}}} K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_i, \bar{\mathbf{y}}))$ and $\epsilon_i = \min_{\bar{\mathbf{y}}} l(\bar{\mathbf{y}}, \mathbf{y}_i) - \sum_{j, \tilde{\mathbf{y}}} \mu_{j, \tilde{\mathbf{y}}} K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_j, \tilde{\mathbf{y}}))$. Step 3 in Algorithm 1 decrease the dual objective function at least by $\frac{\epsilon_i^2}{2R_i}$.*

PROOF. Let $\mu_{i, \bar{\mathbf{y}}}$ be one of the variables which are moved in Step 5. By Equations (8) and (13), we have $\sum_{j, \tilde{\mathbf{y}}} \mu_{j, \tilde{\mathbf{y}}} K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_j, \tilde{\mathbf{y}})) < l(\bar{\mathbf{y}}, \mathbf{y}_i)$. Lemma 3.3 shows the following decrease of the objective function when optimizing over $\mu_{i, \bar{\mathbf{y}}}$.

$$\frac{(l(\bar{\mathbf{y}}, \mathbf{y}_i) - \sum_{j, \tilde{\mathbf{y}}} \mu_{j, \tilde{\mathbf{y}}} K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_j, \tilde{\mathbf{y}})))^2}{2K((\mathbf{x}_i, \bar{\mathbf{y}}), (\mathbf{x}_i, \bar{\mathbf{y}}))} \geq \frac{\epsilon_i^2}{2R_i}.$$

Optimizing over all the variables in μ_B in Step 3 further decreases the dual objective function. \square

Since the dual objective function is convex and quadratic, and the feasible solution region is bounded, the dual objective function is also bounded. Therefore, it is expected that the size of μ_B is bounded according to Theorem 3.4, as shown in the following theorem.

THEOREM 3.5. *Let $R = \max_i R_i$, $\epsilon = \min_i \epsilon_i$, and $\Delta = \max_{i, \bar{\mathbf{y}}} l(\bar{\mathbf{y}}, \mathbf{y}_i)$. There are at most $2nRC\Delta\epsilon^{-2}$ constraints in set μ_B when Algorithm 1 terminates.*

PROOF. The minimum of the primal objective function in Equation (5) can be upper bounded as follows:

$$\min \left(\frac{1}{2} \|\alpha\|^2 + C \sum_{i=1}^n \xi_i \right) \leq \frac{1}{2} \|\alpha\|^2 + C \sum_{i=1}^n \max_{\bar{\mathbf{y}}} (l(\bar{\mathbf{y}}, \mathbf{y}_i) - \alpha^\top \Phi_{i, \mathbf{y}_i, \bar{\mathbf{y}}})$$

$$\leq C \sum_{i=1}^n \max_{\bar{\mathbf{y}}} l(\bar{\mathbf{y}}, \mathbf{y}_i) \leq nC\Delta.$$

According to the relations between the dual problem and primal problem, the minimum of the dual objective function is lower bounded by $-nC\Delta$. When set μ_B is empty, the value of the dual objective function is zero. Therefore, Algorithm 1 terminates after incrementally adding at most $2nRC\Delta\epsilon^{-2}$ constraints according to Theorem 3.4. \square

Note that the number of the constraints in set μ_B does not depend on the original number of the constraints in the primal problem, which is a crucial property of the EML algorithm since the number of the constraints in the primal problem could be increased combinatorially or infinite in many real-world problems.

3.5. Comparison with Other Methods

In the max margin optimization problem Equation (5), only some of the constraints determine the optimal solution. We call these constraints active constraints. Other constraints are automatically met as long as these active constraints are valid. The EMML algorithm uses this fact to solve the optimization problem by substantially reducing the number of the dual variables in Equation (7).

In the recent literature, there are also other methods attempting to reduce the number of the constraints. Taskar et al. [2005] reduce the number of the constraints by considering the dual of the loss-augmented problem. However, the number of the constraints in their approach is still intractably large for a large structured output space and a large training set. They do not use the fact that only some of the constraints are active in the optimization problem.

Tsochantaridis et al. [2004] propose a cutting plane algorithm that finds a small set of active constraints. One issue of this algorithm is that it needs to compute the most violated constraint for each training instance, which would involve another optimization problem in the output space. In EMML, instead of selecting the most violated constraint, we arbitrarily select a constraint that violates the optimality condition of the optimization problem. Thus, the selection of the constraint does not involve any optimization problem. The SVM's computation complexity $R(k)$ is between k^2 and k^3 , where k is the number of constraints. EMML and Tsochantaridis's method could find different set of active constraints, but it is reasonable to assume that the number of constraints T from these two different sets are roughly the same for the same optimization problem. From Equation (3), the number of constraints for one training instance is C_r^m , where m is the number of unique words, r is the average number of the annotation words for one image, and C_r^m is the number of r -combinations from m unique words. Therefore, the computation complexity for the EMML is $\sum_{i=1}^T R(i)$. In Tsochantaridis et al. [2004], Tsochantaridis first selects the most violated constraint for each training instance, which involves the SVM optimization problem with C_r^m constraints. So, the computation complexity to select the most violated constraints is $nhR(C_r^m)$, where n is the number of training instances and h is the number of iterations. So, Tsochantaridis's method has the total computation complexity of $nhR(C_r^m) + \sum_{i=1}^T R(i)$. When $r = \frac{m}{2}$, C_r^m increases exponentially with m . Therefore, EMML is much more efficient in learning with a much faster convergence rate than Tsochantaridis et al. [2004].

4. MULTIMODAL DATA MINING

The solution to the Lagrange dual problem makes it possible to capture the semantic relationships among different data modalities. We show that the developed EMML framework can be used to solve for the general multimodal data mining problem in all the scenarios. Specifically, given a training dataset, we immediately obtain the direct relationship between the VRep space and the word space using the EMML framework in Algorithm 1. Given this obtained direct relationship, we show below that all the multimodal data mining scenarios concerned in this article can be facilitated.

4.1. Image Annotation

Image annotation refers to generating annotation words for a given image. First, we partition the input image into blocks and compute the feature vector in the feature space for each block. We then compute the similarity between feature vectors and the VReps in terms of the distance. We return the top n most-relevant VReps. For each VRep, we compute the score between this VRep and each word as the function f in Equation (1). Thus, for each of the top n most relevant VReps, we have a ranking-list of the words in terms of the score. We then merge these n ranking lists and sort them

to obtain the overall ranking list of the whole word space. Finally, we return the top m words as the annotation result. In our experiments, m and n are fixed constants. Although different images have different numbers of annotation words in real-world image databases, people are always interested in the most relevant annotation results. Thus, small, fixed m and n are appropriate for the problem considered in this article.

In this approach, the score between the VReps and the words is computed in advance. Thus, the computation complexity of image annotation is only related to the number of the VReps. Under the assumption that all the images in the image database follow the same distribution, the number of the VReps is independent of the database scale. Therefore, the computation complexity in this approach is $O(1)$, which is independent of the database scale.

4.2. Word Query

Word query refers to generating corresponding images in response to a query word. For a given word input, we compute the score between each VRep and the word as the function f in Equation (1). Thus, we return the top n most relevant VReps. Since for each VRep, we compute the similarity between this VRep and each image in the image database in terms of the distance, for each of those top n most relevant VReps, we have a ranking list of images in terms of the distance. Then, we merge these n ranking lists and sort them to obtain the overall ranking-list in the image space. Finally, we return the top m images as the query result.

For each VRep, the similarity between this VRep and each image in the image database is computed in advance. Similar to the analysis in Section 4.1, the computation complexity is only related to the number of the VReps, which is $O(1)$.

4.3. Image Retrieval

Image retrieval refers to generating semantically similar images to a query image. Given a query image, we annotate it using the procedure in Section 4.1. In the image database, for each annotation word j , there are a subset of images S_j in which this annotation word appears. We then have the union set $S = \cup_j S_j$ for all the annotation words of the query image.

On the other hand, for each annotation word j of the query image, the word query procedure in Section 4.2 is used to obtain the related sorted image subset T_j from the image database. We then merge these subsets T_j to form the sorted image set T in terms of their scores. The final image retrieval result is $R = S \cap T$.

In this approach, the synergy between the image space and the word space is exploited to reduce the semantic gap based on the developed learning approach. Since the complexity of the retrieval methods in Sections 4.1 and 4.2 are both $O(1)$, and since these retrievals are only returned for the top few items, respectively, finding the intersection or the union is $O(1)$. Consequently, the overall complexity is also $O(1)$.

4.4. Multimodal Image Retrieval

The general scenario of multimodal image retrieval is a query as a combination of a series of images and a series of words. Clearly, this retrieval is simply a linear combination of the retrievals in Sections 4.2 and 4.3 by merging the retrievals together based on their corresponding scores. Since each individual retrieval is $O(1)$, the overall retrieval is also $O(1)$.

4.5. Inference Between Different Stages

For a fruit fly embryo image database such as the Berkeley Drosophila embryo image database that is used for part of our experimental evaluations, we have embryo images classified in advance into different stages of the embryo development with separate sets

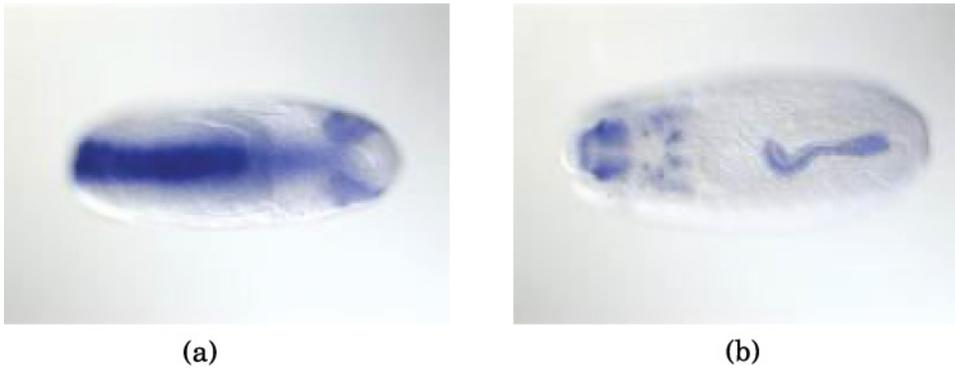


Fig. 2. A pair of embryo images corresponding to the same gene in the two different stages.

of textual words as the annotation to those images in each of these stages. In general, images in different stages may or may not have the direct semantic correspondence (e.g., they all correspond to the same gene), not even speaking that images in different stages may necessarily exhibit any visual similarity. Figure 2 shows an example of a pair of embryo images at stages 9 and 10 [Figure 2(a)] and stages 13–16 [Figure 2(b)], respectively. They both correspond to the same gene in the two different stages.¹ However, it is clear that they exhibit a very large visual dissimilarity.

Consequently, it is not appropriate to use any pure visual feature based similarity retrieval method to identify such image-to-image correspondence between different stages. Furthermore, we also expect to have the word-to-image and image-to-word inference capabilities across different stages, in addition to the image-to-image inference.

Given this consideration, this is exactly where the proposed approach for multimodal data mining can be applied to complement the existing pure retrieval-based methods to identify such correspondence. Typically in such a fruit fly embryo image database, there are textual words for the annotation to the images in each stage. These annotation words in one stage may or may not have the direct semantic correspondence to the images in another stage. However, since the data in all the stages are from the same fruit fly embryo image database, the textual annotation words between two different stages share a semantic relationship that can be obtained by a domain ontology.

In order to apply our approach to this inference problem between different stages, we treat each stage as a separate multimedia database, and map the inference problem to a retrieval-based multimodal data mining problem by applying the approach to the two stages such that we take the multimodal query as the data from one stage and pose the query to the data in the other stage for the retrieval-based multimodal data mining. Figure 3 illustrates the diagram of the two stages (stage i and stage j , where $i \neq j$) image-to-image inference.

Clearly, in comparison with the retrieval-based multimodal data mining analyzed in the previous sections, the only additional complexity here in the inference between different stages is the inference part using the domain ontology in the word space. Typically this ontology is small in scale. In fact, in our evaluations for the Berkeley *Drosophila* embryo image database, this ontology is handcrafted and is implemented as a look-up table for word matching through an efficient hashing function. Thus, this part of the computation may be ignored. Consequently, the complexity of the inference

¹The Berkeley *Drosophila* embryo image database is given in such a way that images from several real stages are mixed together to be considered as one “stage.” Thus, stages 9 and 10 are considered as one stage, and so are stages 13–16.

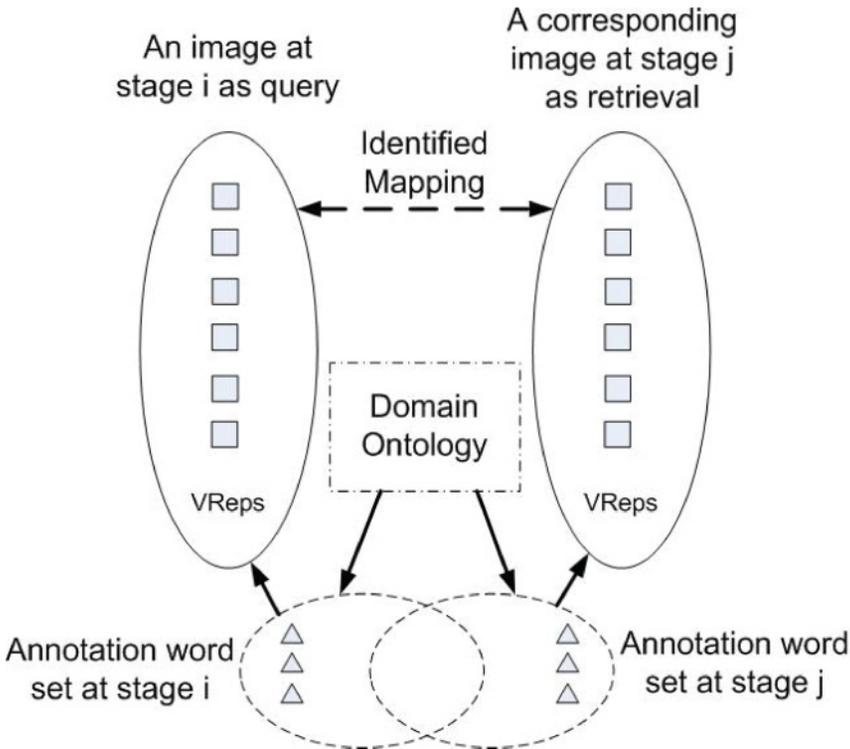


Fig. 3. An illustrative diagram for image-to-image across two stages inference.

between different stages-based multimodal data mining is the same as that of the retrieval-based multimodal data mining, which is independent of database scale.

5. EMPIRICAL EVALUATIONS

While EMML is a general learning framework, and it can also be applied to any structured prediction problem, for the evaluation purpose, we apply it to the multimodal data mining from image databases. In recent literature, Corel image database has been extensively used to evaluate the image retrieval performance [Duygulu et al. 2002; Barnard et al. 2003; Feng et al. 2004]. It has been argued [Westerveld and de Vries 2003] that the Corel data are relatively easy to annotate and retrieve due to its small number of concepts and small variations of visual contents. In order to truly capture the difficulties in real scenarios, we apply EMML to two different real-world image databases. Specifically, we investigate the multimodal retrieval on a real-world image database (called WEB database hereafter), which is from various crawled Web pages. The multimodal data mining including the multimodal retrieval and inference between different stages is evaluated on the Berkeley Drosophila embryo image database.

5.1. Multimodal Retrieval on WEB Database

WEB database consists of a collection of real-world images automatically crawled from the Web. The images and the surrounding text describing the image contents in the Web pages are extracted from the blocks containing the images by using the VIPS algorithm [Cai et al. 2003]. The surrounding text is processed by the standard text processing techniques to obtain the annotation words. We compare EMML framework with a state-of-the-art multimodal data mining method MBRM [Feng et al. 2004] for the retrieval performance.

Table I. Examples of the Image Annotation Produced by EMML and MBRM on WEB Database

| Method |  |  |  |
|--------|---|---|--|
| EMML | *green select american send *picture | zoom *color bush bird *butterfly | *bird zoom send bed ecard |
| MBRM | zoom spacer ecard send bed | zoom button spacer ecard send | spacer ecard send zoom *bird |

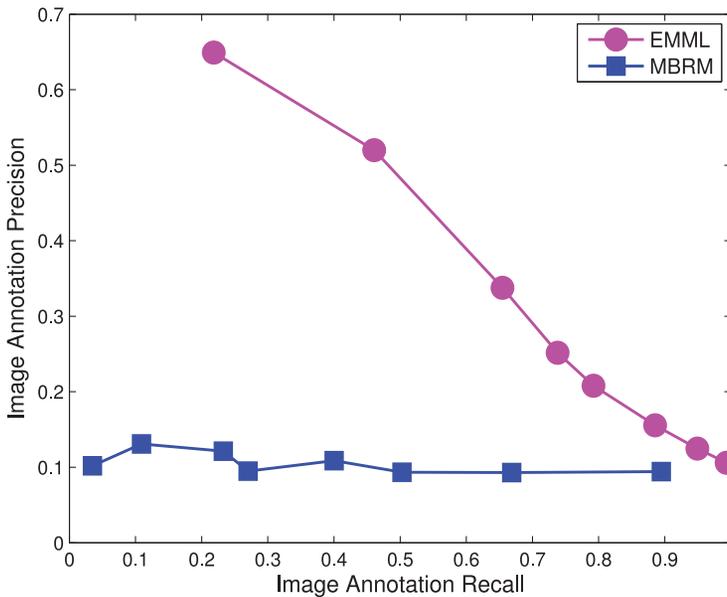


Fig. 4. Precisions and recalls of image annotation between EMML and MBRM on WEB database.

WEB database consists of 5,084 images and 753 annotation words. However, most of the annotation words only occur in a few images. We remove these annotation words and also these images are removed if their annotation becomes empty. Then, a collection of 4,402 images and 51 annotation words is obtained after this preprocessing step.

We then split the whole database into two parts (one-third and two-thirds), with the two-thirds used in the training and the one-third used in the evaluation testing. The EMML framework is applied to obtain the model parameters.

We investigate the performance from top 2 to top 50 retrieval results. In the figures below, the horizontal axis denotes the recall and the vertical axis denotes the corresponding precision for the top retrieval results. Table I shows some examples of the image annotation obtained by EMML framework and the MBRM model on the test image set. Top five words are taken as the annotation of the image. We use the symbol * to denote the relevant word that occurs in the ground truth of the image annotation. Figure 4 reports the precisions and recalls averaged over 1,320 queries for

Table II. Examples of the Single Word Query Produced by EMML and MBRM on WEB Database

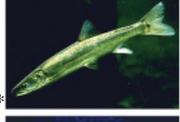
| Method | bird | fish | snake |
|--------|---|---|--|
| EMML |  * |  |  |
| |  |  |  * |
| |  * |  |  |
| |  |  |  * |
| |  |  |  |
| MBRM |  * |  |  |
| |  |  |  * |
| |  * |  * |  * |
| |  * |  * |  * |
| |  * |  * |  * |

image annotation in comparison with the MBRM model. Table II shows some examples of the word query obtained by EMML framework and the MBRM model, where top five images are taken as the query result of the word and the symbol * is used to denote the relevant image that has the word in the ground truth of its annotation. Figure 5 reports the precisions and recalls averaged over 51 queries for word query in comparison with MBRM model. Table III shows some examples of the image retrieval obtained by EMML framework and the MBRM model on the test image set, where top five images

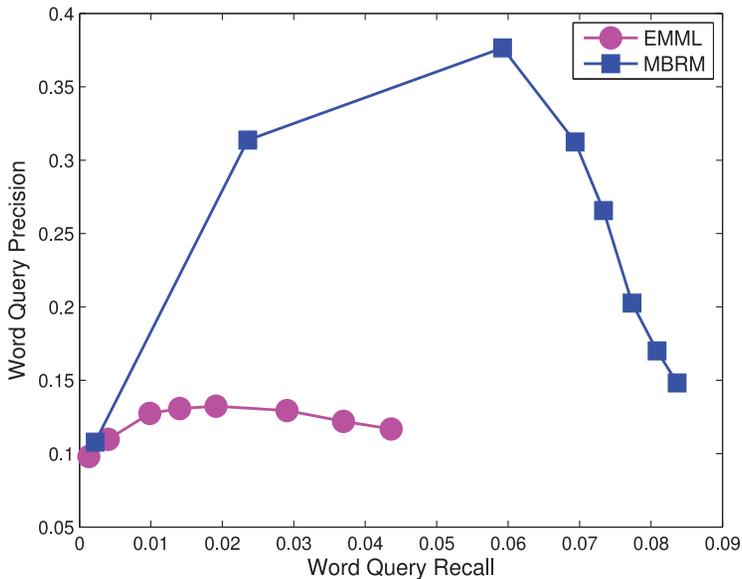


Fig. 5. Precisions and recalls of single word query between EMML and MBRM on WEB database.

are taken as the retrieval result of the image and the symbol * is used to denote the relevant image which has the common annotation words with the input image in the ground truth. Figure 6 reports the precisions and recalls averaged over 1,320 queries for image retrieval in comparison with the MBRM model.

In summary, there is no single winner for all the cases. Overall, EMML outperforms MBRM substantially in the scenarios of image annotation and image retrieval. In the scenario of single word query, MBRM performs much better than EMML. One possible reason is that the number of annotation words is limited compared to the number of images, and EMML is not able to fully capture the relations among images and words. To apply EMML to larger datasets including more annotation words and images is one of our future research directions.

5.2. Multimodal Data Mining on Drosophila Embryo Image Database

The Berkeley Drosophila embryo image database² is used for the multimodal data mining task defined in this article. We evaluate EMML's performance using this database for both the retrieval-based and the inference-based multimodal data mining scenarios. As in Section 5.1, we compare EMML framework with the MBRM model for the mining performance.

In this image database, there are total 16 stages of the embryo images archived in six different folders with each folder containing two to four real stages of the images; there are total 36,628 images and 227 words in all the six folders; not all the images have annotation words. For the retrieval-based multimodal data mining evaluations, we use the fifth folder as the multimedia database, which corresponds to stages 11 and 12. There are about 5,500 images that have annotation words and there are 64 annotation words in this folder. We split the whole folder's images into two parts (one-third and two-thirds), with the two-thirds used in the training and the one-third used in the evaluation testing. For the inference-based multimodal data mining evaluations, we

²<http://www.fruitfly.org>.

Table III. Examples of the Image Retrieval Produced by EMML and MBRM on WEB Database

| | | | |
|--------|---|---|--|
| Method |  |  |  |
| EMML |  *  *  *  *  * |  *  *  *  *  * |  *  *  *  *  * |
| MBRM |      |  *  *    |    *   |

use the fourth and the fifth folders for the two stages inference evaluations, and use the third, the fourth, and the fifth folders for the three stages inference evaluations. Consequently, each folder here is considered as a “stage” in the inference-based multimodal data mining evaluations. In each of the inference scenarios, we use the same

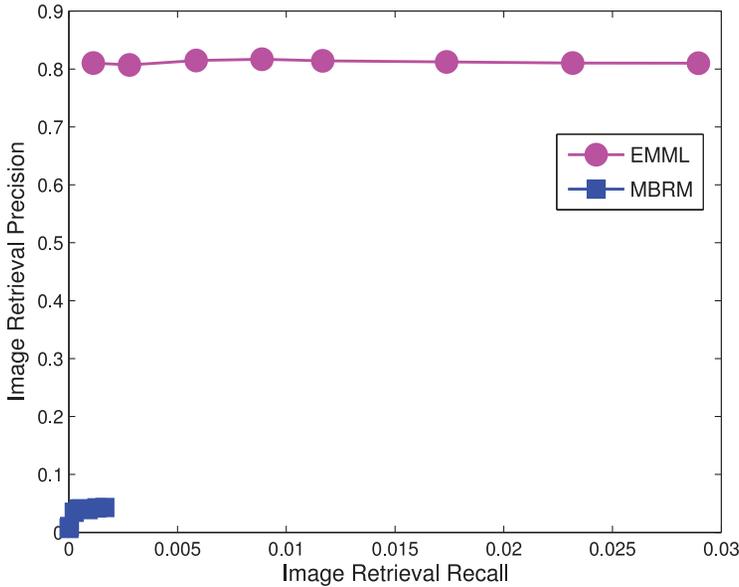


Fig. 6. Precisions and recalls of image retrieval between EMML and MBRM on WEB database.

split as we do in the retrieval-based multimodal data mining evaluations for training and testing.

In order to facilitate the inference capabilities, we handcraft the ontology of the words involved in the evaluations. This is simply implemented as a simple look-up table indexed by an efficient hashing function. For example, *cardiac mesoderm primordium* in the fourth folder is considered as the same as *circulatory system* in the fifth folder. With this simple ontology and word matching, the proposed approach may be well applied to this inference problem between different stages for the multimodal data mining.

As in Section 5.1, the EMML algorithm is applied to obtain the model parameters. Figure 7 reports the precisions and recalls averaged over 1,648 queries for image annotation in comparison with MBRM model. Table IV shows some examples of the image annotation obtained by EMML framework and the MBRM model on the test image set. Top five words are taken as the annotation of the image. We use the symbol * to denote the relevant word that occurs in the ground truth of the image annotation. Similarly, Figure 8 reports the precisions and recalls averaged over 64 queries for word query in comparison with the MBRM model. Table V shows some examples of the word query obtained by EMML framework and the MBRM model, where top five images are taken as the query result of the word and the symbol * is used to denote the relevant image that has the word in the ground truth of its annotation. Figure 9 reports the precisions and recalls averaged over 1,648 queries for image retrieval in comparison with the MBRM model. Table VI shows some examples of the image retrieval obtained by EMML framework and the MBRM model on the test image set, where top five images are taken as the retrieval result of the image and the symbol * is used to denote the relevant image that has the common annotation words with the input image in the ground truth. In these tables, PR stands for primordium and SA stands for specific anlage. As shown in these Figures and Tables, EMML demonstrates significant performance improvement over MBRM. For example, in image retrieval (Figure 9 and Table VI), EMML is able to successfully retrieve the semantically relevant images

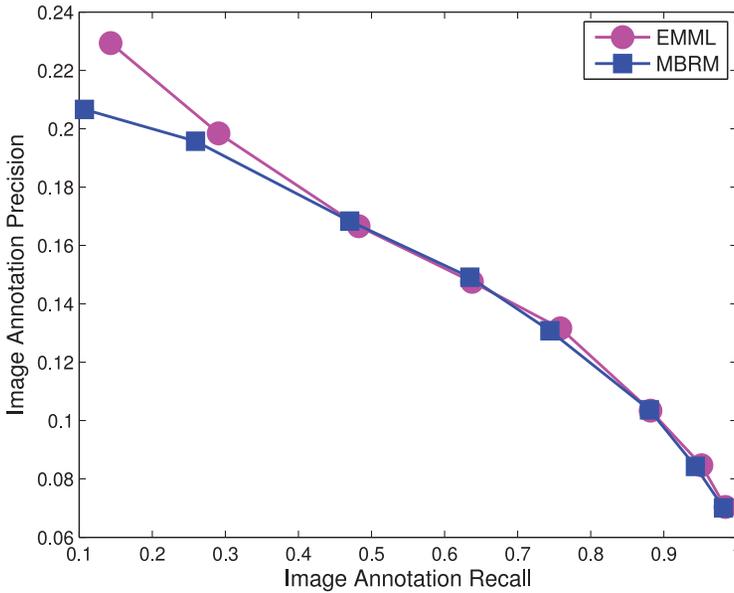


Fig. 7. Precisions and recalls of image annotation between EMML and MBRM on Drosophila Embryo Image Database.

Table IV. Examples of the Image Annotation Produced by EMML and MBRM on Drosophila Embryo Image Database

| Method |  |  |  |
|--------|--|--|---|
| EMML | posterior midgut PR *trunk mesoderm PR anterior midgut PR *lateral cord neuron *protocerebrum PR | *posterior midgut PR trunk mesoderm PR *anterior midgut PR amnioserosa lateral cord neuron | *posterior midgut PR *trunk mesoderm PR *anterior midgut PR *protocerebrum PR lateral cord neuron |
| MBRM | posterior midgut PR *trunk mesoderm PR anterior midgut PR *lateral cord neuron hindgut proper PR | *posterior midgut PR trunk mesoderm PR *anterior midgut PR lateral cord neuron hindgut proper PR | *posterior midgut PR *trunk mesoderm PR *anterior midgut PR lateral cord neuron hindgut proper PR |

even if the images are very different from each other visually, whereas MBRM fails to generate the semantically similar images to a query image, because EMML has the better capability to capture the semantic relationship among different data modalities.

For the two-stage inference, Figure 10 reports the precisions and recalls averaged over 1,648 queries for image-to-word inference in comparison with the MBRM model, and Figure 11 reports the precisions and recalls averaged over 64 queries for word-to-image inference in comparison with the MBRM model. Figure 12 reports the precisions and recalls averaged over 1,648 queries for image-to-image inference in comparison with the MBRM model. Finally, for the three-stage inference, Figure 13 reports precisions and recalls averaged over 1,100 queries for image-to-image inference in comparison with the MBRM model.

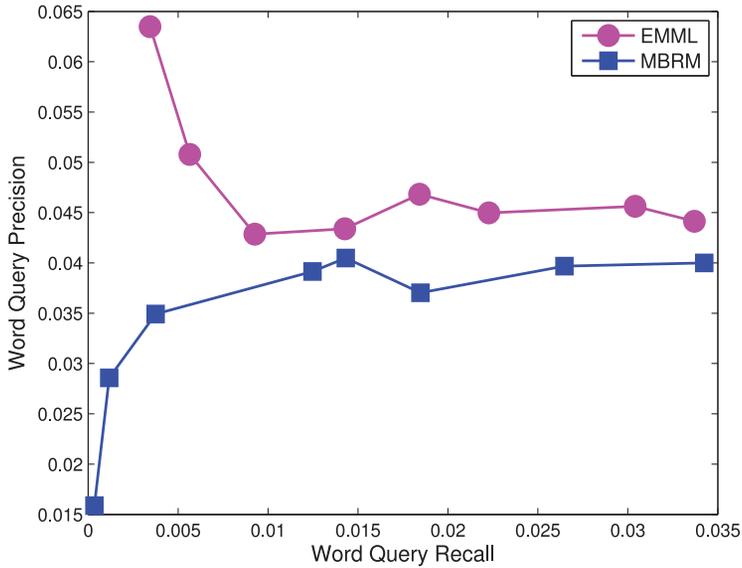
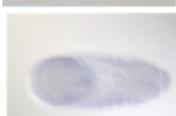
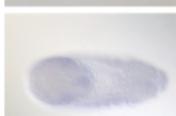
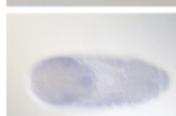


Fig. 8. Precisions and recalls of word query between EMLL and MBRM on Drosophila Embryo Image Database.

As in Section 5.1, there is no single winner for all the cases. Overall, EMLL outperforms MBRM substantially in the scenarios of word query and image retrieval, and slightly in the scenario of two-stage word-to-image inference and three-stage image-to-image inference. On the other hand, MBRM has a slight better performance than EMLL in the scenario of two-stage image-to-word inference. For all other scenarios, the two methods have a comparable performance. For the inference between different stages, domain ontology is very important in order to obtain the good performance and the superior domain ontology requires advanced domain knowledge from experts. Since we are not the experts on the Drosophila Embryo, the constructed domain ontology might have some defects, which could lead to the worse performance of EMLL than MBRM in the case of two-stage image-to-world inference.

In the Table IV, MBRM gives the same annotation words in the same order for three different images, which is due to the underlying generative process in the MBRM model. MBRM assumes that images are independent of each other and words are also independent of each other. This assumption leads to an issue in computing the posterior probability of the words conditioned on the given image. According to $p(w|x) \propto p(x|w)p(w)$, one expects that a popular word (larger $p(w)$) is likely to have a larger posterior probability because $p(x|w)$ is independent of the popularity of the word. The Drosophila Embryo database is unbalanced in the sense that some words occur more frequently than others. Therefore, MBRM gives the larger posterior probability to the more popular words. Similarly for the Table V, MBRM gives the same images in the same order for three different words. The reason behind this is the image property of the Drosophila Embryo database. All of the embryo images are in the light color, so the feature vectors computed from images are overwhelmed by the light color information. Due to the same reason in the image annotation, those images containing mostly light color will have high posterior probability. Since EMLL considers the multimodal data mining problem as the structured predication problem where the images and words are related to each other, this posterior probability issue does not exist. The examples in the Tables I–III

Table V. Examples of the Single Word Query Produced by EMMML and MBRM on Drosophila Embryo Image Database

| Method | dorsal epidermis PR | posterior spiracle SA | trunk mesoderm PR |
|--------|---|---|--|
| EMML |  |  |  |
| |  |  |  |
| |  |  |  |
| |  |  |  |
| |  |  |  |
| MBRM |  |  |  |
| |  |  |  |
| |  |  |  |
| |  |  |  |
| |  |  |  |

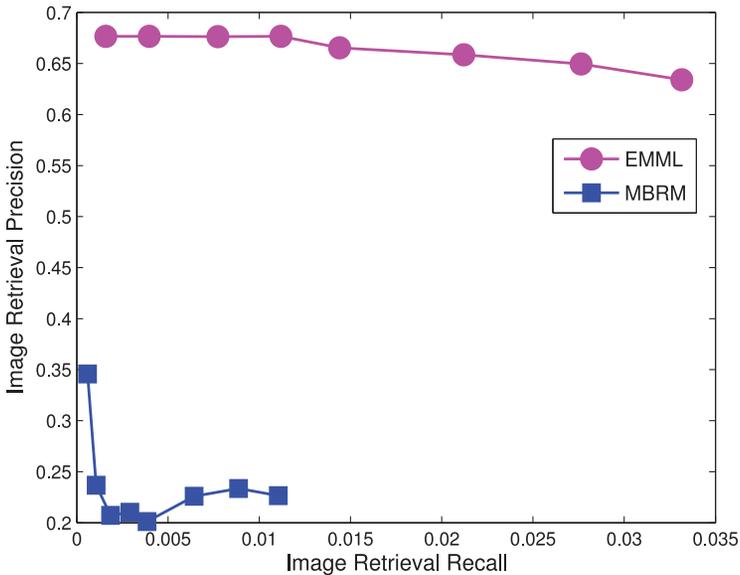


Fig. 9. Precisions and recalls of image retrieval between EMML and MBRM on Drosophila Embryo Image Database.

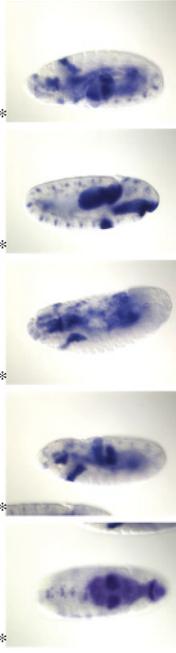
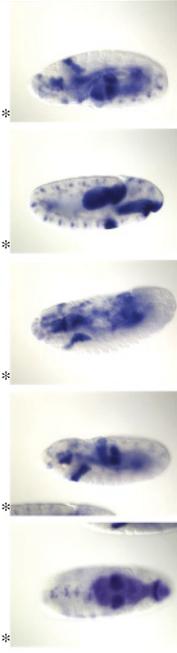
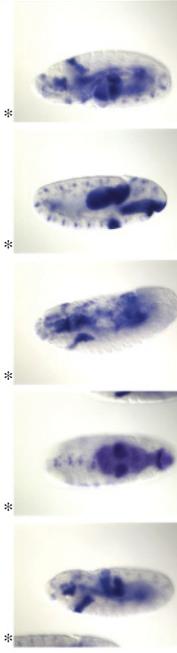
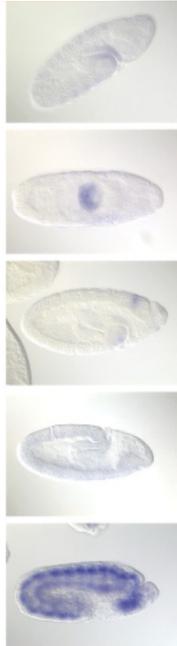
do not have this problem because the words and images' color in the WEB database are more diversified.

To investigate that whether EMML improves the performance over MBRM or not from the viewpoint of statistics, we perform the paired hypothesis tests based on the results. Since EMML and MBRM have a comparable performance in the two-stage inference and three-stage inference, we only perform the paired hypothesis tests for the image annotation, word query, and image retrieval. Two hypothesis tests are performed: paired right-tail t-test and paired two-sided Wilcoxon signed-rank test, where the null hypothesis is that the difference between the results of the two methods comes from a distribution with zero mean and the alternative hypothesis is that the mean is greater than zero (right-tail t-test) or is not zero (signed-rank test). According to the p -value show in the Table VII, the null hypothesis are rejected for most of the scenarios, which indicates that EMML statistically improves the performance by modeling the multimodal data mining problem as a structured prediction problem.

In order to demonstrate the strong scalability of EMML approach to multimodal data mining, we take image annotation as a case study and compare the scalability between EMML and MBRM. We randomly select three subsets of the embryo image database in different scales (1,000, 2,000, 3,000 images, respectively), and apply both methods to the subsets to measure the query response time. The query response time is obtained by taking the average response time over 1,648 queries. Since EMML is implemented in MATLAB environment and MBRM is implemented in C in Linux environment, to ensure a fair comparison, we report the scalability as the relative ratio of a response time to the baseline response time for the respective methods. Here, the baseline response time is the response time to the smallest scale subset (i.e., 1,000 images). Table VIII documents the scalability comparison. Clearly, MBRM exhibits a linear scalability w.r.t the database size, while that of EMML is constant. This is consistent with the scalability analysis in Section 4.

In order to verify the fast learning advantage of EMML in comparison with the existing max margin-based learning literature, we have implemented one of the most

Table VI. Examples of the Image Retrieval Produced by EMLL and MBRM on Drosophila Embryo Image Database

| Method |  |  |  |
|--------|---|---|--|
| EMML |  |  |  |
| MBRM |  |  |  |

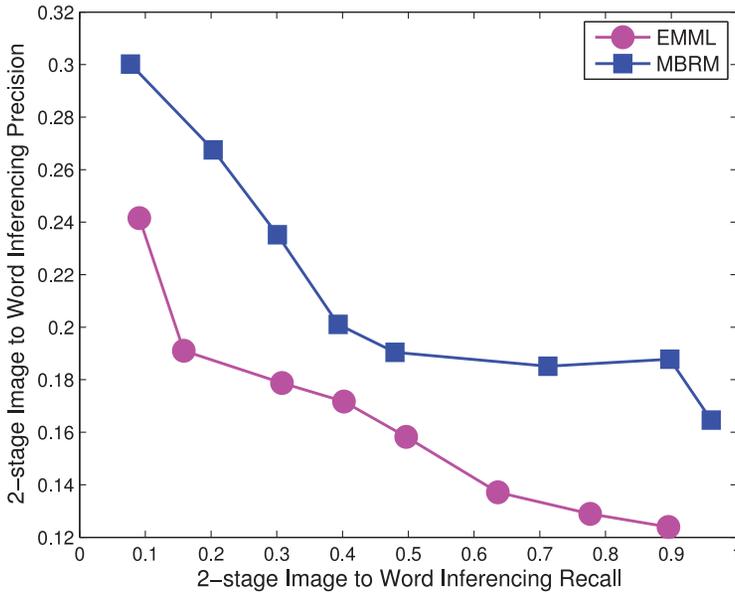


Fig. 10. Precisions and recalls of two-stage image to word inference between EMML and MBRM on Drosophila Embryo Image Database.

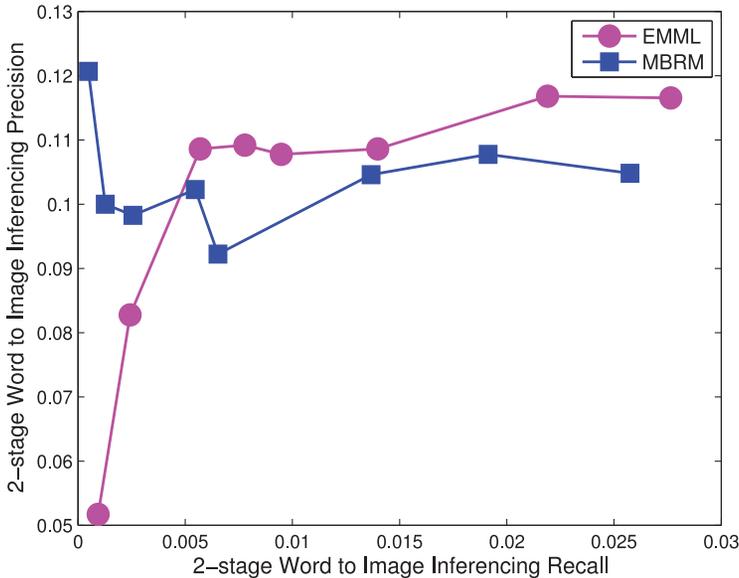


Fig. 11. Precisions and recalls of two-stage word to image inference between EMML and MBRM on Drosophila Embryo Image Database.

recently proposed max margin learning methods by Taskar et al. [2005]. For the reference purpose, in this article, we call this method as TCKG. We have applied both EMML and TCKG to a small dataset randomly selected from the whole Berkeley embryo database, consisting of 110 images along with their annotation words. The reason we use this small dataset for the comparison is that we have found that in MATLAB platform TCKG immediately runs out of memory when the dataset is larger, due to

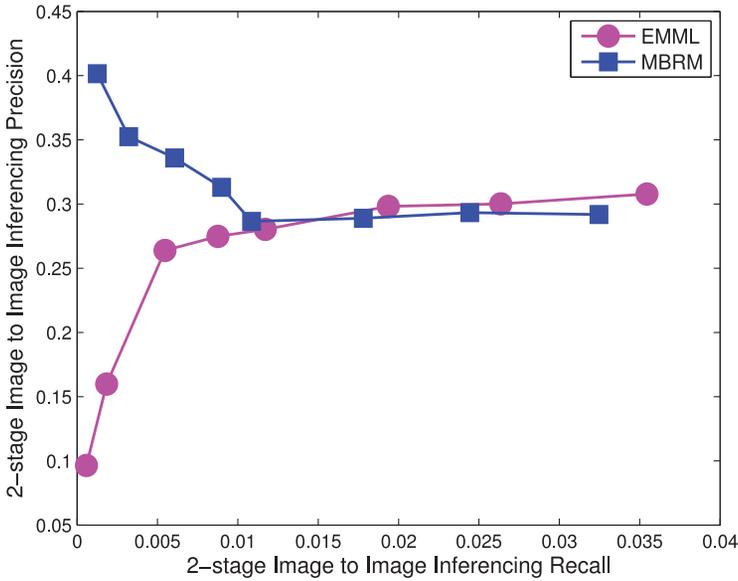


Fig. 12. Precisions and recalls of two-stage image to image inference between EMML and MBRM on Drosophila Embryo Image Database.

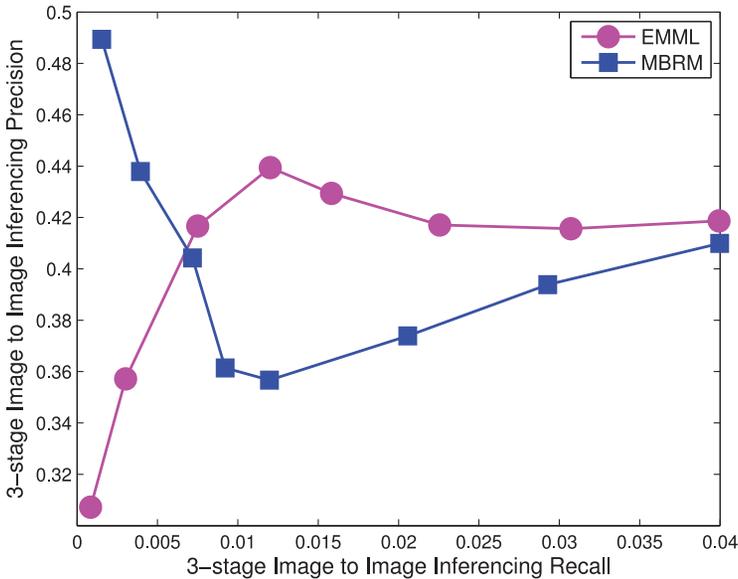


Fig. 13. Precisions and recalls of three-stage image to image inference between EMML and MBRM on Drosophila Embryo Image Database.

the large number of the constraints, which is typical for the existing max margin learning methods. Under the environment of 2.2-GHz CPU and 1-GB memory, TCKG takes about 14hours to complete the learning for such a small dataset while EMML only takes about 10minutes. We have examined the number of the constraints reduced in both methods during their executions for this dataset. EMML has reduced the number of the constraints in a factor of 70 times more than that reduced by TCKG.

Table VII. p -Value with the Significance Level 0.05

| | | paired t-test | | signed-rank test | |
|-----------------|------------------|---------------|--------|------------------|--------|
| | | Precision | Recall | Precision | Recall |
| Embryo Database | Image Annotation | 0.0107 | 0.0078 | 0.1730 | 0.0107 |
| | Word Query | 0.0009 | 0.0156 | 0.0205 | 0.0009 |
| | Image Retrieval | 0.0027 | 0.0078 | 0.0001 | 0.0027 |
| Web Database | Image Annotation | 0.0078 | 0.0078 | 0.0127 | 0.0001 |
| | Word Query | 0.0078 | 0.0078 | 0.9962 | 0.9967 |
| | Image Retrieval | 0.0078 | 0.0078 | 0.0001 | 0.0045 |

Table VIII. Comparison of Scalability

| Database Size | 1000 | 2000 | 3000 |
|---------------|------|------|------|
| EMML | 1 | 1 | 1 |
| MBRM | 1 | 2.1 | 3.2 |

This explains why EMML is about 70 times faster than TCKG in learning for this dataset.

6. CONCLUSION

We have developed a new max margin learning framework—the EMML, and applied it to developing an effective and efficient multimodal data mining solution. EMML attempts to find a small set of active constraints, and thus, is more efficient in learning than the existing max margin learning literature. Consequently, it has a much faster convergence rate that is verified in empirical evaluations. The multimodal data mining solution based on EMML is highly scalable in the sense that the query response time is independent of the database scale. This advantage is also supported through the complexity analysis as well as empirical evaluations. While EMML is a general learning framework and can be used for general structured predication problem, for the evaluation purpose, we have applied it to two different real-world image databases and have reported the evaluations against a state-of-the-art multimodal data mining method.

REFERENCES

- Yasemin Altun, Ioannis Tsochantaridis, and Thomas Hofmann. 2003. Hidden Markov support vector machines. In *Proceedings ICML*. Washington, DC.
- Kobus Barnard, Pinar Duygulu, David Forsyth, Nando de Freitas, David M. Blei, and Michael I. Jordan. 2003. Matching words and pictures. *J. Mach. Learn. Res.* 3, 3 (2003), 1107–1135.
- D. Blei and M. Jordan. 2003. Modeling annotated data. In *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM New York, NY, USA, 127–134.
- Stephen Boyd and Lieven Vandenberghe. 2004. *Convex Optimization*. Cambridge University Press.
- Ulf Brefeld and Tobias Scheffer. 2006. Semi-supervised learning for structured output variables. In *Proceedings ICML*. Pittsburgh, PA.
- D. Cai, S. Yu, J.-R. Wen, and W.-Y. Ma. 2003. VIPS: A vision-based page segmentation algorithm. Microsoft Technical Report MSR-TR-2003-79.
- E. Chang, Kingshy Goh, G. Sychay, and Gang Wu. 2003. CBSA: content-based soft annotation for multimodal image retrieval using Bayes point machines. *IEEE Trans. Circuits and Systems for Video Technology* 13, 1 (January 2003), 26–38.
- W. Chu, Z. Ghahramani, and D. L. Wild. 2004. A graphical model for protein secondary structure prediction. In *Proceedings ICML*. Banff, Canada.
- R. Datta, W. Ge, J. Li, and J. Z. Wang. 2006. Toward bridging the annotation-retrieval gap in image search by a generative modeling approach. In *Proceedings ACM Multimedia*. Santa Barbara, CA.

- Hal Daume III and Daniel Marcu. 2005. Learning as search optimization: Approximate large margin methods for structured prediction. In *Proc. ICML*. Bonn, Germany.
- Pinar Duygulu, Kobus Barnard, Nando de Freitas, and David Forsyth. 2002. Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In *Proceedings 7th European Conference on Computer Vision*, Vol. IV. 97–112.
- S. L. Feng, R. Manmatha, and V. Lavrenko. 2004. Multiple Bernoulli relevance models for image and video annotation. In *Proceedings International Conference on Computer Vision and Pattern Recognition*. Washington, DC.
- Yoav Freund and Robert E. Schapire. 1999. Large margin classification using the perceptron algorithm. In *Mach. Learn.* 37, 3 (1999), 277–296.
- Zhen Guo, Zhongfei (Mark) Zhang, Eric P. Xing, and Christos Faloutsos. 2007a. Enhanced max margin learning on multimodal data mining in a multimedia database. In *Proceedings the 13th ACM International Conference on Knowledge Discovery and Data Mining*. San Jose, CA.
- Zhen Guo, Zhongfei (Mark) Zhang, Eric P. Xing, and Christos Faloutsos. 2007b. A max margin framework on image annotation and multimodal image retrieval. In *Proceedings the IEEE Annual International Conference on Multimedia and Expo*. Beijing, China.
- John Lafferty, Andrew McCallum, and Fernando Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings ICML*. 282–289.
- Andrew McCallum, Dayne Freitag, and Fernando Pereira. 2000. Maximum Entropy Markov models for information extraction and segmentation. In *Proceedings ICML*. 591–598.
- Edgar Osuna, Robert Freund, and Federico Girosi. 1997. An improved training algorithm for support vector machines. In *Proceedings of IEEE NNSP'97*. Amelia Island, FL.
- J.-Y. Pan, H.-J. Yang, C. Faloutsos, and P. Duygulu. 2004. Automatic multimedia cross-modal correlation discovery. In *Proceedings of the 10th ACM SIGKDD Conference*. Seattle, WA.
- A. W. M. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain. 2000. Content-based image retrieval at the end of the early years. *IEEE Trans. Pattern Anal. Mach. Intell.* 22, 12 (2000), 1349–1380.
- B. Taskar, C. Guestrin, and D. Koller. 2003. Max-margin Markov networks. In *Proceedings Neural Information Processing Systems Conference*. Vancouver, Canada.
- B. Taskar, V. Chatalbashev, D. Koller, and C. Guestrin. 2005. Learning structured prediction models: A large margin approach. In *Proceedings ICML*. Bonn, Germany.
- Ioannis Tsochantaridis, Thomas Hofmann, Thorsten Joachims, and Yasemin Altun. 2004. Support vector machine learning for interdependent and structured output spaces. In *Proceedings ICML*. Banff, Canada.
- Vladimir Naumovich Vapnik. 1995. *The Nature of Statistical Learning Theory*. Springer.
- T. Westerveld and A. de Vries. 2003. Experimental evaluation of a generative probabilistic image retrieval model on ‘easy’ data. In *Proceedings of the SIGIR Multimedia Information Retrieval Workshop 2003*.
- Yi Wu, Edward Y. Chang, and Belle L. Tseng. 2005. Multimodal metadata fusion using causal strength. In *Proceedings ACM Multimedia*. Hilton, Singapore, 872–881.

Received September 2009; revised December 2014; accepted February 2015