

## Exploring Interaction Between Images and Texts for Web Image Categorization

Lei Li<sup>1\*</sup>, Wenting Lu<sup>2†</sup>, Jingxuan Li<sup>1</sup>, Tao Li<sup>1</sup>

<sup>1</sup>School of Computing and Information Sciences  
Florida International University  
Miami, Florida 33199, United States

Honggang Zhang<sup>2</sup>, Jun Guo<sup>2</sup>

<sup>2</sup>Pattern Recognition and Intelligence System Lab  
Beijing University of Posts and Telecommunications  
Beijing 100876, P.R.China

### Abstract

With the rapid development of technologies for fast access to the Internet and the popularization of digital cameras, enormous digital images are posted and shared online everyday. Simultaneously, web images are usually organized by topics of events and are often assigned appropriate topic-related text descriptions. Given a set of images along with corresponding texts, a challenging problem is how to utilize the available information to perform image retrieval tasks, such as image classification and image clustering. Previous works on image categorization focus on either adopting text or image features, or simply combining these two types of information together. In this paper, we propose two novel approaches (**Dynamic Weighting** and **Region-based Semantic Concept Integration**) to categorize the images under the “supervision” of topic-related text descriptions; In addition, we provide a comparative experimental investigation on utilizing text and image information to tackle image classification. Empirical experiments on a manually collected image dataset (consisting of images related to the events after disasters) demonstrate the efficacy of our proposed classification methods.

### Introduction

Multimedia information plays an increasingly important role in human’s daily activities. With the rapid development of technologies of fast access to the Internet and the popularization of digital cameras, enormous digital images are posted and shared online everyday. Besides great convenience, how to retrieve images that satisfy the needs of web users in multimedia databases is becoming more and more difficult and challenging. Particular, web image categorization, as a crucial step of image retrieval, attracts much more attention and is very useful in the subsequent procedures, such as indexing and organizing web image databases, browsing and searching web images, and discovering interesting patterns from images (Yin et al. 2009).

However, web image categorization is not a trivial task due to the diversity of image content and the limited information. In general, the images posted on the Internet may have great visual differences, rendering image categoriza-

tion challenging since it is difficult to extract common features shared by most images as the comparison base. Fortunately, text information is often provided by web users to describe the general contents of images, *e.g.*, image titles, headers, or text descriptions assigned to them. A possible solution to web image categorization is to take such text information into consideration. Specifically, we can initially extract different types of features (image and text features), and then categorize the images by delving into the special characteristics of the integrated version of different features.

In this paper, we explore the feasibility of using text information as a “guidance” for image categorization by proposing two novel methods (**Dynamic Weighting** and **Region-based Semantic Concept Integration**), which achieve better performance comparing with existing approaches. Specifically, the proposed *Dynamic Weighting* assumes that different image features might carry semantic meanings with different significance, and under the “supervision” of text features, the importance of different image features can be dynamically decided. It is straightforward that the important features may have more distinguishable power for image categorization. Another proposed method *Region-based Semantic Concept Integration* first segments images into different regions, and then categorizes images based on the correlation between regions and semantic concepts. Moreover, we provide a comparative experimental study on integrating text and image information to perform image categorization. Empirical experiments on a manually collected image dataset (including images related to the events after disasters) demonstrate the efficacy of our proposed methods.

The rest of this paper is organized as follows. In Section 2 we review some related works that combine image and text features to perform image classification tasks. In Section 3 we give algorithmic details of the two proposed approaches to effectively integrate text information with image information for image categorization. Section 4 presents a detailed experimental comparison among different approaches and finally we conclude the paper in Section 5.

### Related Work

Most of the existing web image categorization approaches often focus on utilizing text descriptions of images to categorize images via simply matching keywords. Currently, the majority of web search engines still adopt this technique due

\*Li, L. and Lu, W. contributed equally to this work.

†This author is a visiting Ph.D student.

to its fast speed. There are some limitations in text-based image categorization: (1) web images cannot be appropriately classified if there is no text information assigned to them; (2) the manual text labeling is too subjective due to human assignments, which might result in bias or noise to web image categorization; (3) using a few words to describe the content of an image is not enough since the limited text description can only provide a relatively sparse feature space. Therefore, the performance of traditional text-based web image categorization systems is very limited.

To solve the above problems, many research publications (Blei and Jordan 2003; Giacinto, Roli, and Fumerga 2002; Kalva, Enembreck, and Koerich 2007; Wu et al. 2004; Zhu, Yeh, and Cheng 2006) aim to design multi-view learning algorithms to learn classifiers from multiple information sources via integrating different types of features together to perform classification. In general, multi-view learning methods can be categorized into three different groups:

1. **Feature Integration:** Enlarge the feature representation to incorporate all attributes from different sources and produce a unified feature space. The advantage of feature integration is that the unified feature representation is often more informative and also allows many different data mining methods to be applied and systematically compared. One disadvantage is the increased learning complexity and difficulty as the data dimension becomes large (Wu, Oviatt, and Cohen 2002);
2. **Semantic Integration:** Keep data intact in their original form and computational methods are applied to each feature space separately. Results on different feature spaces are then combined by either voting (Carter, Dubchak, and Holbrook 2001), Bayesian averaging (Bishop 2006), or the hierarchical expert system approach (Jordan and Jacobs 1994). One advantage of semantic integration is that it can implicitly learn the correlation structure between different sets of features (Li and Ogihara 2005).
3. **Intermediate Integration:** A compromise between the feature integration and the semantic integration. The idea is to keep the feature spaces in their original forms and integrate them at the similarity computation or the Kernel level (Schölkopf and Smola 2002; Lanckriet et al. 2004). Different weights can be assigned different data sources. Standard computational methods can then be applied once the total similarity is computed.

**Our contribution:** In order to explore the feasibility of using texts as a guidance for image classification, we propose two novel multi-view learning methods to achieve better performance for web image categorization. Further, we present an empirical investigation on different methods for combining text information with image features and compare their classification performance.

## Classification Algorithms

Web image categorization is a key step for many web-based multimedia applications. It is crucial for the subsequent processes (e.g., image retrieval) and has a direct impact on the speed and accuracy of other applications related to images

on the web. As mentioned above, three different multi-view learning approaches have been used to resolve the problem of web image categorization. However, all these approaches focus on simply combining two data sources (text and image information), and none of them takes the advantage that one data source can provide “guidance” for another on how to perform categorization task. In this paper, besides providing a comparative study of the previous works, we also propose two novel web image categorization methods – *Dynamic Weighting* and *Region-based Semantic Concept Integration* – which employ the text-based information (*i.e.*, the text itself and the semantic concepts hidden in the text) to guide the classification, and consequently achieve better image categorization results.

## Dynamic Weighting

Image feature extraction techniques tend to extract a huge number of image features based on different criteria. Among these features, some of them might carry significant semantic information about the image, whereas some others might be less crucial for the tasks being executed on the image. Particularly in image classification, the extracted features should be more representative and carry substantial amount of semantic meanings. Therefore, it might be helpful to dynamically assign different weights to different image features so that the features with more importance can be captured and play more meaningful roles on the classification. Some previous works (Shao et al. 2009) on music information retrieval demonstrate how to learn appropriate similarity metrics based on the correlation between acoustic features and user access patterns. Motivated by this, we incorporate the concept of dynamic feature weighting into our image classification problem.

Specifically in image classification, given that human perception of an image is well approximated by its text description, a good weighting schema for the extracted image features guided by text information may lead to a good similarity measurement, and therefore better classification results. Let  $\mathbf{m}_i = (\mathbf{f}_i, \mathbf{t}_i)$  denote the  $i$ -th image in the image collection, where  $\mathbf{f}_i$  and  $\mathbf{t}_i$  represent its image features and text features respectively. Let  $S_f(\mathbf{f}_i, \mathbf{f}_j; \mathbf{w}) = \sum_l f_{i,l} f_{j,l} w_l$  be the image-based similarity measurement between the  $i$ -th and the  $j$ -th images when the parameterized weights are given by  $\mathbf{w}$ . Let  $S_t(\mathbf{t}_i, \mathbf{t}_j) = \sum_k t_{i,k} t_{j,k}$  be the similarity measurement between the  $i$ -th and the  $j$ -th images based on their text description features, in general, the words with specific meanings extracted from texts. Here for each  $k$ ,  $t_{i,k}$  denotes whether the  $k$ -th word appears in the text description of the  $i$ -th image. To learn appropriate weights  $\mathbf{w}$  for image features, we can enforce the consistency between similarity measurements  $S_f(\mathbf{f}_i, \mathbf{f}_j; \mathbf{w})$  and  $S_t(\mathbf{t}_i, \mathbf{t}_j)$ . The above idea leads to the following optimization problem:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i \neq j} (S_f(\mathbf{f}_i, \mathbf{f}_j; \mathbf{w}) - S_t(\mathbf{t}_i, \mathbf{t}_j))^2 \quad s.t. \mathbf{w} \geq 0. \quad (1)$$

Let  $p$  be the number of image features. The summation in

Eq.(1) can be rewritten as follows:

$$\begin{aligned}
& \sum_{i \neq j} (S_f(\mathbf{f}_i, \mathbf{f}_j; \mathbf{w}) - S_t(\mathbf{t}_i, \mathbf{t}_j))^2 \\
&= \sum_{i \neq j} \left( f_{i,1}f_{j,1}w_1 + \dots + f_{i,p}f_{j,p}w_p - \sum_k t_{i,k}t_{j,k} \right)^2 \\
&= \sum_{i \neq j} \left( (f_{i,1}f_{j,1}w_1 + \dots + f_{i,p}f_{j,p}w_p)^2 \right. \\
&\quad \left. - 2(f_{i,1}f_{j,1}w_1 + \dots + f_{i,p}f_{j,p}w_p) \right. \\
&\quad \left. \times \left( \sum_k t_{i,k}t_{j,k} \right) + \left( \sum_k t_{i,k}t_{j,k} \right)^2 \right),
\end{aligned}$$

where  $f_{i,l}$  is the  $l$ -th feature in the image feature set  $f_i$  and  $f_{j,l}$  is the  $l$ -th feature in the image feature set  $f_j$ . Let  $n$  be the number of images, and let

$$F = \begin{bmatrix} f_{1,1}f_{2,1} & f_{1,2}f_{2,2} & \dots & f_{1,g}f_{2,g} \\ \dots & \dots & \dots & \dots \\ f_{n-1,1}f_{n,1} & f_{n-1,2}f_{n,2} & \dots & f_{n-1,g}f_{n,g} \end{bmatrix},$$

and

$$T = \begin{bmatrix} i=j f_{i,1}f_{j,1} \binom{n}{k} t_{i,k}t_{j,k} \\ \vdots \\ i=j f_{i,g}f_{j,g} \binom{n}{k} t_{i,k}t_{j,k} \end{bmatrix},$$

where  $F$  is a  $\binom{n}{2} \times p$  matrix and  $T$  is a  $p \times 1$  matrix. Thus, Eq.(1) is equivalent to

$$\begin{aligned}
\mathbf{w}^* &= \operatorname{argmin} \left[ \frac{1}{2} \times 2(Fw)^T(Fw) - T^T w \right] \\
&= \operatorname{argmin} \left[ \frac{1}{2} (w^T(2F^T F)w + (-2T^T)w) \right] \text{ s.t. } \mathbf{w} \geq 0.
\end{aligned}$$

This optimization problem can be addressed using quadratic programming techniques (Gill, Murray, and Wright 1981). After calculating the dynamic weights for each image features, we multiply the feature values with the corresponding weights for each image, and finally obtain a new feature space with weighted information. These feature vectors can then be fed into classifiers. In the experiments, we will show how much the classification results are influenced by our dynamic weighting schema.

## Region-based Semantic Concept Integration

In the real-world applications, an image always contains various semantic concepts and these concepts often intersect with each other, which is not helpful to efficiently extract semantic information. In this section, we explore the feasibility of utilizing the underlying semantic concepts of text information as a ‘‘guidance’’ to facilitate image categorization. To address the issue mentioned above, we firstly divide original images into different regions to ensure that the content of each region represents almost the same local pattern, and then based on the local semantic patterns of the images, we propose our *Region-based Semantic Concept Integration* method. Figure 1 shows the framework of our proposed approach, which can be divided into four different

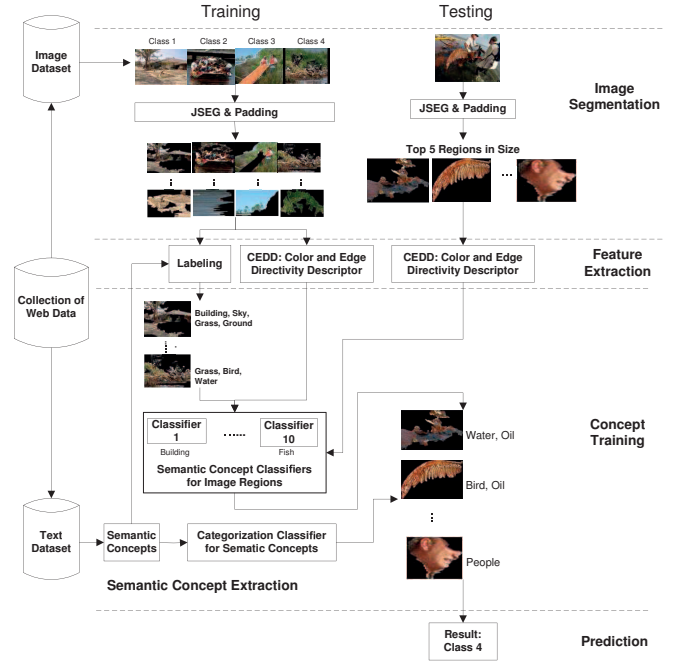


Figure 1: Framework of region-based concept integration.

sub-processes: *semantic concept extraction*, *image segmentation*, *feature extraction* on each region, and *region-based semantic concept classification*. In the following, we provide algorithmic details of these four processes.

**Semantic Concept Extraction** In this process, we initially analyze the text description related to each image, and then obtain some original high-frequency terms in these texts by using MALLET (McCallum 2002), a java-based package for statistical natural language processing. We then compare the semantics of these high-frequency words and summarize them to generate several most general semantic concepts using WordNet (Miller 1995). The general concepts are represented as the hypernyms of the high-frequency words. Then each text description can be represented by the combination of these concepts. The generalized concepts can provide guidance on how to select image region samples in the training step of semantic concept classifier, as well as to train a concept model that builds the relationship between semantic concepts and original categories, which will be described in Section ‘‘Design of Semantic Concept Classifiers’’.

**Image Segmentation** In order to associate the images with the generalized concepts extracted from the procedure of ‘‘Semantic Concept Extraction’’, we need to segment images into different regions such that each region can be related to one or more semantic concepts. Ideally, image segmentation aims to divide original web images into different regions based on the criterion that each region contains only one object or one part of an object. However, due to the limitation of current segmentation techniques, it is very difficult to perfectly segment images (Liu, Zhang, and Lu 2008;

Li, Socher, and Fei-Fei 2009). In this paper, we use a state-of-the-art segmentation method JSEG (Joint Systems Engineering Group) (Deng and others 2001), which segments images based on color and texture information. In the segmentation algorithm, image color space is first quantized into several classes, and a color class-map of each image is then obtained via re-representing each pixel of the original image by its corresponding color class label. After that, the spatial segmentation is performed on this color class-map which can be viewed as a special type of texture composition. Here, a criterion named “ $J$ -Value” is used to measure whether the segmentation is reasonable or not. If one image consists of several homogeneous color regions, the color classes will be separated from each other and the value of  $J$  is larger. Figure 2 presents some examples of the segmented images of web images using JSEG.



Figure 2: An illustration of image segmentation.

**Region Feature Extraction** After segmenting images into different regions, image feature extraction is performed on each region. Note that color and texture are two of the most general global features in the field of image processing and computer vision, both of which have their own advantages and drawbacks. In order to represent image/region effectively, we adopt CEDD (Color and Edge Directivity Descriptor) (Chatzichristofis and Boutalis 2008), which is a new low-level feature descriptor incorporating both color feature and texture feature. In CEDD, a novel but effective method is adopted to integrate a 24-bins color histogram and a 6-bins texture histogram to form a final 144-bins histogram. One of the most important characteristics of CEDD is its low computational power needed for feature extraction, in comparison with the needs of most of MPEG-7 descriptors. For detailed algorithmic procedure of CEDD, please refer to (Chatzichristofis and Boutalis 2008).

**Design of Semantic Concept Classifiers** In this step, the design of semantic concept classifiers can be divided into the following two parts:

*Semantic Concept Classifiers for Image Regions:* Based on previous steps of image segmentation and feature extraction, we obtain a set of regions and their corresponding 144-dimension feature vectors. Then we need to train  $N$  (the number of generalized semantic concepts obtained from text descriptions) different semantic concept classifiers respectively using training regions. Each of the training regions is

manually labeled with multiple concepts, but each semantic concept classifier is just designed for binary classification identifying whether a region contains this concept or not. Therefore, when a testing region arrives, it will be fed into these  $N$  classifiers respectively and a  $N$ -dimension vector is obtained for the input of the categorization classifier.

*Categorization Classifier for Semantic Concept:* We design a multi-label classifier to build the relationship between semantic concepts and the original categories. Here we use One-Against-One (OAO) (Hsu and Lin 2002) method to design the multi-label SVM classifier. OAO designs an original binary SVM classifier between two random classes of samples, and it needs  $k(k-1)/2$  original binary SVM classifiers. In addition, we use a voting strategy in classification, in which each binary classification is considered to be a voter where votes can be cast for all the regions, and finally each region is designated to be in a class with the maximum number of votes.

When a testing image comes, it will be segmented into several regions. Then the system will choose the top  $M$  regions in size automatically (if the actual number of regions for the image is less than  $M$ , the system will adopt the actual number of regions) and extract CEDD feature from these regions. Note that the larger the region is, the more information it contains, and if the region is too small, it is difficult to identify the exact semantic meaning and therefore may involve noises. After that, each region will pass through the  $N$  different semantic concept classifiers respectively, to identify whether this region contains the concept or not. If the region contains this concept, a concept label will be assigned to this region so that each region will be assigned multiple concept labels, and the region can be represented by a  $N$ -dimension vector. Then we integrate these  $M$   $N$ -dimension vectors into a single  $N$ -dimension vector describing whether the original image contains certain concepts or not. Finally, we use the categorization classifier for semantic concepts to predict which class the original testing image belongs to.

## Experiments

### Real World DataSet

Unlike the traditional scene object databases, which are mainly focusing on visual categorization, web images are usually organized by topics of events. From a classification perspective, the differences between these two kinds of images are as follows:

- Images in the same category for visual scene categorization are visually similar, but object scaling, rotating, occluding and submerging often happen in clutter background; Comparatively, web images in the same category may vary visually but very similar in terms of semantic concepts;
- Images in the same visual category contain the same object or scene and it would occupy most area of the whole image, whereas web image focus on reflecting just one aspect of the whole event;

Therefore in our work, we focus on the conceptual information contained in images, but not simple rotation or zooming on the same objects.

**Dataset Description** To systematically study the previous multi-modal feature combination methods and compare them with our proposed approaches, we manually collected 355 colored web images and their corresponding text descriptions about “the aftermath of disasters”, which include 4 different topics: Hurricane building collapse, Hurricane flood, Oil spill seagrass, and Oil spill animal death. Each category includes 101, 101, 53 and 100 images respectively, and the entire image dataset is split into two parts: about 70% images (247 in total) are randomly selected for training and the rest 30% images (108 in total) are taken as the test data.

Note that for *Region-based Semantic Concept Integration* method, we generalize 10 (*i.e.*,  $N=10$ ) concepts from the text description set, including “building”, “water”, “sky”, “grass”, “oil”, “bird”, “ground”, “people”, “helicopter” and “fish”. In addition, we generate 1573 regions from 247 training images with the guidance of these semantic concepts, and 513 regions from 108 testing images by automatically choosing the top 5 (*i.e.*,  $M=5$ ) regions in size (if the actual number of regions is less than 5, the system will adopt the actual number of regions).

## Design of Experiments

In our experiment, we use LIBSVM (Chang and Lin 2001) as our base classification tool. The parameter tuning is done via  $k$ -fold cross validation. For the purpose of comparison, we first implement five existing methods for web image categorization, then compare their classification performance with those of our proposed approaches. These five existing methods include:

- Text-based Classification (*Text* for short): Extract text features from texts assigned to the corresponding images, and then use these features to feed SVM classifier;
- Image-based Classification (*Img* for short): Extract CEDD image features from the images, and then use these features to feed SVM classifier;
- Feature Integration (*Feat* for short): Treat unique terms as text features and extract image features using CEDD. Note that the extracted CEDD feature is a 144-dimension vector, while the cardinality of the text features is 1788. To balance the contribution of different features to the classification results, we choose the top 144 terms with high frequency as the text features. We combine the features of text and image together by simply concatenating these two types of features to form a 288-dimension vector as the input of SVM classifier.
- Semantic Integration (*Sem* for short): Train two classifiers based on text features and image features (CEDD) respectively, and then ensemble these two classifiers to be an integrated version, similar to the method proposed in (Carter, Dubchak, and Holbrook 2001).
- Intermediate Integration (*Sim* for short): Compute the pairwise similarity using text-based features and image-based features (CEDD) respectively, and then use the weighted summation of these two types of similarities as the similarity measurement between images. Note that

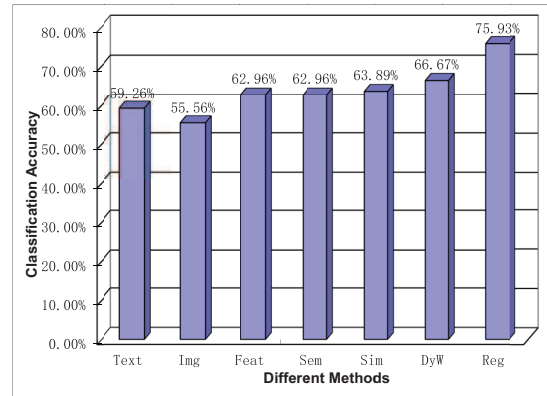


Figure 3: Comparison among the five existing methods and our two proposed methods based on the accuracy of web image categorization.

different weights can be assigned to the features of different data sources. We tune the weight factor to find the optimal one through empirical comparison.

## Experiment Results

**Comparison among all the methods** In Figure 3, the comparison among the above 5 different methods and our proposed methods (*DyW* and *Reg* for short) based on the accuracy of web image categorization results are presented.

From the comparison results, we observe that the best performance of web image categorization using single-modal approaches is less than 60%. However, once the text and image data sources are integrated using two-modal information fusion techniques, the categorization performance is improved. The intuitive explanation for the improvement is that two-modal approaches are able to incorporate the advantages of the two data sources together, which leads to better categorization results compared with using only one type of features (here the “advantage” represents the positive contribution of the features of certain data source to the image categorization results). In addition, compared with the categorization results provided by the five existing methods, our proposed approaches outperform others. The reason behind the performance improvement is straightforward: our two proposed methods share one common characteristic – employing advantages of one data source to enrich the other data source. In other words, these two methods explore inherent connections between two data sources by utilizing text-based information to either find out the best weighting schema for the image-based features or generate the classifiers from the semantic image regions related to text concepts.

**Comparison between *DyW* and *Reg*** We further compare two proposed approaches based on their classification performance on each category of images. The comparison results are shown in Table 1. From the results, we have the following observations:

	Dynamic Weighting			Region-based Concept		
	Precision	Recall	F1-measure	Precision	Recall	F1-measure
Category 1	0.4565	0.6774	0.5455	0.6429	0.8710	0.7397
Category 2	0.6875	0.7097	0.6984	0.7391	0.5484	0.6296
Category 3	1.0000	0.0625	0.1176	1.0000	0.5000	0.6667
Category 4	0.9655	0.9333	0.9492	0.8571	1.0000	0.9231

Table 1: Comparison of classification results on each category using our proposed methods. Note that “Category 1-4” represents Hurricane building collapse, Hurricane flood, Oil spill seagrass, and Oil spill animal death respectively.

1. *DyW* and *Reg* provide reasonable performance on Category 1 and 2, and good results on Category 4.
2. The performance of these two methods on Category 1 and 2 is slightly worse than the results on Category 4. This is due to the characteristics of our image dataset. Most of the images in Category 1 and 2 contain a lot of “semantic noise”. For instance, most of the images in Category 1 focus on the collapse “buildings”, but “grass” and “water” also appear in these images. These “noise” would cause our classification methods to misclassify these images into the other categories. Even though *Reg* incorporates the semantic information and shows better results than *DyW*, “noise” still exists to some extent.
3. The recall of both two methods on Category 3 is very low. After analysis, we found that most of the images in Category 3 are about “grass”; however, “grass” appears in almost all the categories, which results in the misclassification of the images.
4. *Reg* outperforms *DyW* on category 1 and 3, and the performance of *Reg* on category 2 and 4 is comparable with *DyW*. The reason that the overall performance of *Reg* is better than the one of *DyW* is that *Reg* could benefit from the semantic information hidden in the text whereas *DyW* only make use of raw text information.

## Conclusions

In this paper, we study the problem of combining two data sources (text and image) to perform image categorization tasks and show that such combination can lead to better classification results comparing with using individual data sources. Also, we propose two novel multi-view learning methods which can effectively utilize the image-related text data to find out better schemas to classify the images. The empirical results show that our proposed methods outperform the previous methods in terms of the accuracy of classification results, and they can provide solid basis for the subsequent procedures of image retrieval.

## Acknowledgement

This work is partially supported by the U.S. Department of Homeland Security’s VACCINE Center under Award Number 2009-ST-061-CI0001 and by the Army Research Office under grant number W911NF-10-1-0366, and by the Fundamental Research Funds for the Central Universities, 111 Project of China (B08004) and Scientific Research Foundation for the Returned Overseas Chinese Scholars, State Edu-

cation Ministry. We would also like to thank the Distributed Multimedia Information Systems Laboratory at Florida International University for helping us collect the dataset.

## References

- Bishop, C. 2006. *Pattern recognition and machine learning*.
- Blei, D., and Jordan, M. 2003. Modeling annotated data. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval*, 127–134. ACM.
- Carter, R.; Dubchak, I.; and Holbrook, S. 2001. A computational approach to identify genes for functional RNAs in genomic sequences. *Nucleic Acids Research* 29(19):3928.
- Chang, C., and Lin, C. 2001. LIBSVM: a library for support vector machines.
- Chatzichristofis, S., and Boutalis, Y. 2008. Cedd: Color and edge directivity descriptor: A compact descriptor for image indexing and retrieval. In *Proceedings of the 6th International Conference on Computer Vision Systems*, 312–322.
- Deng, Y., et al. 2001. Unsupervised segmentation of color-texture regions in images and video. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 800–810.
- Giacinto, G.; Roli, F.; and Fumerga, G. 2002. Unsupervised learning of neural network ensembles for image classification. In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks*, volume 3, 155–159. IEEE.
- Gill, P.; Murray, W.; and Wright, M. 1981. Practical optimization.
- Hsu, C., and Lin, C. 2002. A comparison of methods for multiclass support vector machines. *IEEE Transactions on Neural Networks* 13(2):415–425.
- Jordan, M., and Jacobs, R. 1994. Hierarchical mixtures of experts and the EM algorithm. *Neural computation* 6(2):181–214.
- Kalva, P.; Enembreck, F.; and Koerich, A. 2007. Web image classification based on the fusion of image and text classifiers. In *Proceedings of the 9th International Conference on Document Analysis and Recognition*, 561–568. IEEE Computer Society.
- Lanckriet, G.; Cristianini, N.; Bartlett, P.; Ghaoui, L.; and Jordan, M. 2004. Learning the kernel matrix with semidefinite programming. *The Journal of Machine Learning Research* 5:27–72.
- Li, T., and Ogihara, M. 2005. Semisupervised learning from different information sources. *Knowledge and Information Systems* 7(3):289–309.
- Li, L.; Socher, R.; and Fei-Fei, L. 2009. Towards total scene understanding: Classification, annotation and segmentation in an automatic framework. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2036–2043. IEEE.
- Liu, Y.; Zhang, D.; and Lu, G. 2008. Region-based image retrieval with high-level semantics using decision tree learning. *Pattern Recognition* 41(8):2554–2570.
- McCallum, A. 2002. MALLETT: A Machine Learning for Language Toolkit. <http://mallet.cs.umass.edu>.
- Miller, G. 1995. WordNet: a lexical database for English. *Communications of the ACM* 38(11):39–41.
- Schölkopf, B., and Smola, A. 2002. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT Press.
- Shao, B.; Ogihara, M.; Wang, D.; and Li, T. 2009. Music recommendation based on acoustic features and user access patterns. *IEEE Transactions on Audio, Speech, and Language Processing* 17(8):1602–1611.
- Wu, Y.; Chang, E.; Chang, K.; and Smith, J. 2004. Optimal multimodal fusion for multimedia data analysis. In *Proceedings of the 12th Annual ACM International Conference on Multimedia*, 572–579. ACM.
- Wu, L.; Oviatt, S.; and Cohen, P. 2002. Multimodal integration—a statistical view. *IEEE Transactions on Multimedia* 1(4):334–341.
- Yin, Z.; Li, R.; Mei, Q.; and Han, J. 2009. Exploring social tagging graph for web object classification. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 957–966. ACM.
- Zhu, Q.; Yeh, M.; and Cheng, K. 2006. Multimodal fusion using learned text concepts for image categorization. In *Proceedings of the 14th Annual ACM International Conference on Multimedia*, 211–220. ACM.

# HIERARCHICAL DISASTER IMAGE CLASSIFICATION FOR SITUATION REPORT ENHANCEMENT

Yimin Yang, Hsin-Yu Ha, Fausto Fleites, Shu-Ching Chen, Steven Luis  
School of Computing and Information Sciences  
Florida International University  
Miami, FL 33199, USA  
{yyang010, hha001, fflei001, chens, luiss}@cs.fiu.edu

## Abstract

*In this paper, a hierarchical disaster image classification (HDIC) framework based on multi-source data fusion (MSDF) and multiple correspondence analysis (MCA) is proposed to aid emergency managers in disaster response situations. The HDIC framework classifies images into different disaster categories and sub-categories using a pre-defined semantic hierarchy. In order to effectively fuse different sources (visual and text) of information, a weighting scheme is presented to assign different weights to each data resource depending on the hierarchical structure. The experimental analysis demonstrates that the proposed approach can effectively classify disaster images at each logical layer. In addition, the paper also presents an iPad application developed for situation report management using the proposed HDIC framework.*

## 1. Introduction

Due to the ease of access and wide reach of Internet, more and more multimedia data such as images and videos, along with corresponding textual descriptions, become available through the web everyday. Such availability of content-rich data is extremely valuable for emergency management (EM) personnel as they can take more accurate decisions in disaster situations by having both textual and visual information of the disaster. Nevertheless, currently, EM personnel mostly utilize disaster situation reports (also referred as situation reports) which provide just a textual description of the disaster. To augment situation reports with related disaster images and thus provide EM personnel with images and videos that present valuable information about the disaster, a hierarchical disaster image classification (HDIC) framework is proposed in this paper. Based on multi-source data fusion (MSDF)

and multiple correspondence analysis (MCA) [1], our framework classifies disaster multimedia data into different categories and links these images to related situation reports. In order to obtain images for the disaster domain (i.e., hurricane, oil spill, and earthquake), we collected both images and their corresponding titles and description from Flickr. The HDIC framework utilizes both visual features from images and textual description to demonstrate the performance of combining MCA-based data fusion method with the hierarchical classification approach.

There are two main applications for image classification in the area of disaster analysis: damage detection and damage prediction. Najab [13] used Principal Component Analysis (PCA) to extract the features from remotely-sensed data and classify them into different landcover classes. Gandhe [8] leveraged a framework which includes discrete wavelet transform (DWT) and PCA to help with image mining and weather forecasting, and Hsu [9] applied wavelet transformation, support vector machines, and fuzzy neural networks for image compression, classification and error correction respectively to an intelligent typhoon damage prediction system. In addition, classification of high-resolution disaster images facilitates the process of damage assessment after environmental disasters such as hurricane, tsunami, etc [12, 6, 2, 3, 4]. Unlike the aforementioned works that focus on satellite images [13, 3, 4], images retrieved from multiple remote sensing sensors [12, 6] and aerial photos [9, 2], our framework is able to classify the actual disaster images taken at the disaster location, which have higher complexity and reduce the semantic gap between the images and the disaster categories. In addition, the proposed framework is able to fusion multi-source data in an efficient way achieving higher performance than the individual textual and visual models independently.

The remainder of this paper is organized as follows. Section 2 briefly describes the HDIC framework based on

MSDF and MCA. Section 3 discusses the MCA algorithm for multimedia content analysis. Section 4 presents the details of the visual-text model training. Section 5 discusses the hierarchical classification based on MSDF. Experimental analyses is presented in section 6, and section 7 briefly introduces the ipad application developed based on the HDIC framework. Finally, section 8 concludes the paper.

## 2 HDIC Framework

Depicted in Figure 1, the HDIC framework is composed of two main processes: multi-source model training and hierarchical classification. During the model training process, visual and text features are extracted respectively and fused based on the weighting scheme presented in section 5. Then the models for different categories and sub-categories (subjects) are trained based on the MCA algorithm, generating thresholds for classification. The feature extraction of testing data depends on that of the training data. For example, the discretization intervals of test visual feature should corresponds to those of the training data. Finally, the trained models are applied to the hierarchical classification of images, where the images are firstly classified into general categories, and then passed to the next layer to be assigned to specific subjects.

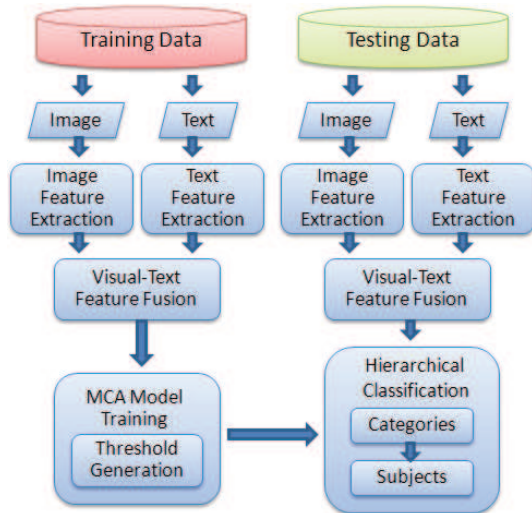


Figure 1: HDIC framework.

## 3 MCA for Multimedia Content Analysis

MCA is an exploratory data analytic technique designed to analyze multi-way tables for some measure of correspondence between the rows and columns [1]. It is a natural extension of the standard correspondence analysis

to more than two variables. The observations used for MCA are a set of nominal variables, each of which is composed of several levels, and each level is coded as a binary value. There is a constraint that one and only one level of the variable gets the value 1. Therefore each observation has the same total, called mass. MCA can also accommodate quantitative variables by recording them as bins, which inspires the idea that it could be applied to numerical data, such as multimedia feature instances. For example, each image feature variable could be discretized into several intervals, and each image can be presented by a series of nominal values.

Motivated by the functionality of MCA as well as its quantitative analysis ability, the utilization of MCA has been explored in our previous works to analyze the data instances described by a set of low-level features to capture the correspondence between items (feature-value pairs) and classes (subjects). The similarity of every item and every class can be presented by the cosine of the angle between each item and class [10, 11]. A smaller angle indicates a higher correlation between the item and class.

## 4 Visual-Text Model Training Based on MCA

This section reveals the feature extraction processes for both visual and text data as well as the model-training procedure based on the MCA algorithm. An iterative threshold determination algorithm is also presented to find out the most appropriate threshold for classification.

### 4.1 Visual feature extraction

There are mainly three steps for visual feature extraction: feature extraction, normalization, and discretization. The first two steps are the same for both training images and test images; however, the discretization of the test images' features is based on the discretized intervals resulted from training image instances.

In order to capture the visual contents of images, two types of feature are extracted: low-level color features and mid-level object location features, which are described as follows:

- *Twelve color features*: black, white, red, red-yellow, yellow, yellow-green, green, green-blue, blue, blue-purple, purple, and purple-red;
- *Nine object location features*: Images are divided into  $3 \times 3$  grid, i.e., nine locations  $L_1, \dots, L_9$ , where  $L_i = 1$  if there is an object whose centroid falls inside  $L_i = 1, 1 \leq i \leq 9$ .

Therefore a total number of 21 features are obtained, where the color features are based on the HSV color space,

and the object location features are extracted using the SPCPE algorithm [5]. Since the color features and object location features are considered equally important, an equal weight (i.e., 0.5) is assigned to each type of features in the normalization step. Finally, an information-gain-based discretization method [7] is used for numerical to nominal transformation.

## 4.2 Text feature extraction

Due to its limitation in descriptive capability, visual features alone could not well represent the content of an image. Therefore text features are introduced to enhance the description. The proposed text feature extraction procedure requires more preprocessing than visual features. First, punctuation characters and stop words are removed, thus obtaining a list of valid words for each image instance. The word frequency is calculated based on all the training instances for each concept (subject). The top N (i.e. 50 in our experiment) words with the highest frequencies are selected as features. A nominal value is assigned to each feature representing the existence or absence of it. Then each image instance could be transformed to a sequence of nominal variables with N dimensions. The feature extraction process of the test data set is almost the same as that of the training data set except for the "get word frequency" step since the construction of testing feature vector is based on the top N words from training data.

## 4.3 Visual-Text Model Training

The process of visual-text model training can be summarized into two major steps: MCA score calculation and threshold generation. More specifically, after visual and text feature extraction of the training data sets, the two sets of feature vectors are concatenated together to form a data set of fused instances, which are used for angle generation based on MCA correlation analysis. The angles, denoted as  $A$ , are calculated using Equation (1), where  $I$  and  $C$  are two-dimensional principal components representing items and classes respectively as described in section 3, and  $j, k$  are indicators of items and features. Then the generated angles are applied to weight conversion as shown in Equation (2). The weight is a measure of the similarity between each item and class. The sum of all of the weights within one instance is denoted as  $S$  (shown in Equation (3)), which is the final evaluation of the relationship between each instance and class. A higher score implies a higher possibility that the instance belongs to the class, which implies the existence of a cut point (threshold) determining the positive or negative attribute of one instance for certain

class (subject).

$$A_k^j = \arccos\left(\frac{I_k^j \cdot C}{|I_k^j| |C|}\right), \quad (1)$$

$$weight_k^j = \pm(1 + \cos(A_k^j \times \pi/180)), \quad (2)$$

$$S_i = \sum_{k=1}^K weight_k^j, i \in \{1, 2, \dots, N\} \quad (3)$$

How to determine the threshold is a critical issue and plays an extremely important role in the final performance of the whole classification algorithm. Therefore an iterative method is designed to find out the threshold for classification based on the training instances:

THRESHOLD-GENERATION:

- 1  $finalF1 = 0;$
- 2  $finalThresh = 0;$
- 3  $sortedScore = \mathbf{sort}(trainScore);$
- 4  $cddThresh = \mathbf{find}(positive);$
- 5 for  $i = 1$  **to length** ( $cddThresh$ )
- 6  $testLabel(1 \mathbf{to} cddThresh(i)) = classLabel;$
- 7  $F1$  calculation;
- 8 **if**  $finalF1 > F1$  ||  $finalF1 - F1 < \gamma$  **then**
- 9  $finalF1 = F1;$
- 10  $finalThresh = sortedScore(cddThresh(i)).$

Steps 1 and 2 initialize the variables of  $finalF1$  and  $finalThresh$ , which store the final F1 score and the corresponding threshold. In step 3, the sort function sorts training scores in descending order, and step 4 finds the indexes of positive scores from the sorted array as candidate thresholds. Step 5 through 10 loop through each candidate to find the best threshold giving the optimal performance in terms of the F1 measure. Specifically, steps 6 and 7 calculate the F1 scores based on precision and recall (refer to section 6). In step 8, the latter condition (i.e.,  $finalF1 - F1$ ) is designed to include the neglected positive instances; it provides the functionality of balancing between recall and precision measures and improves F1 scores. The term  $\gamma$  is a practical parameter, and it is set to be 0.03 in the experiments. Finally, steps 8 and 9 recall the final F1 score and threshold.

## 5 Hierarchical Classification

In order to explore the extensive relationship between various subjects and perform the classification in a more efficient way, a hierarchical classification mechanism is proposed. The hierarchical classification scheme breaks down disaster-related categories into a tree structure which serves to organize general to specific categories. The classification scheme addressed in this paper was developed

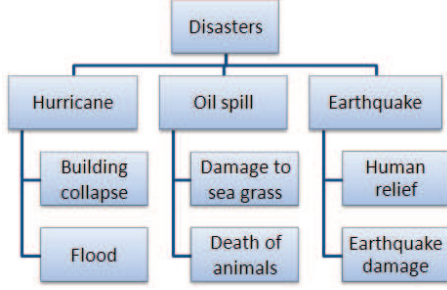


Figure 2: Hierarchical structure.

upon consulting with experts in the disaster management field.

As shown in Figure 2, the top-down category tree classifies images into one of three main categories, i.e., hurricane, oil spill, or earthquake based on text-visual models, and then the chosen category will be further classified into a specific sub-category. Based on the observation that the text data in the second layer has a stronger pattern than that of visual model and vice versa in the third layer, a weighting scheme is proposed to distinguish the significance of visual and text models at different layers and obtain a better fusion result. The fusion score is calculated as follows:

$$score_f = \alpha W_v * score_v + \beta W_t * score_t, \quad (4)$$

$$thresh_f = \alpha W_v * thresh_v + \beta W_t * thresh_t, \quad (5)$$

$$W_v = \frac{F1_v}{F1_v + F1_t}, \quad W_t = \frac{F1_t}{F1_v + F1_t}, \quad (6)$$

$$W_v + W_t = 1, \quad \alpha + \beta = 2. \quad (7)$$

where  $score_v$  and  $score_t$  represent the scores obtained from visual and text models, while  $\alpha W_v$  and  $\beta W_t$  denote the weight factors of visual and text models respectively, and  $score_f$  is the final fused score. The thresholds are fused in the same manner. The  $W_v$  and  $W_t$  are calculated based on the F1 measures of visual and text models at different layers, while the  $\alpha$  and  $\beta$  are tuning parameters. In the experimental analysis, the  $\alpha$  and  $\beta$  are set to be 0.50 and 1.50 in the second layer; 1.23 and 0.77 in the third layer. Finally, the classification rules are generated as follows:

$$finalLabel = \begin{cases} positive, & \text{if } score_f \geq thresh_f \\ negative, & \text{if } score_f < thresh_f \end{cases} \quad (8)$$

## 6 Experimental Analysis

In order to demonstrate the effectiveness of the proposed MCA-based multimedia content analysis, a set of experiments have been conducted to evaluate its

Categories	Subjects	No. of images	Total images
Hurricane (Cat1)	Building collapse (Sub1)	197	1025
	Flood (Sub2)	144	
Oil spill (Cat2)	Damage to sea grass (Sub3)	151	
	Death of animals (Sub4)	176	
Earthquake(Cat3)	Human Relief (Sub5)	221	
	Earthquake damage (Sub6)	136	

Figure 3: Composition of categories and subjects.

performance. The test bed is a web-crawled dataset consisting of 1,025 images with texts downloaded from Flickr. The number of images is limited due to the fact that domain-specific disaster images are not abundant. The images contain three categories and cover six subjects as shown in Figure 3. The categories are denoted as Cat1, Cat2, and Cat3, and the subjects are denoted as Sub1 through Sub6.

In the experimental settings, the hierarchical classification scheme shown in Figure 2 is adopted. Multi-source (text and visual) data fusion is performed at both layer 2 and layer 3. To show the advantages of the multi-source model over single-source models, a comparison between the performances of the multi-source text-visual model and the single-source text and visual models are conducted at each layer. The precision (Equation 9), recall (Equation 10), and F1 (Equation 11) are calculated as the measurements of performance under the 3-fold cross validation approach.

$$precision = \frac{TP}{TP + FP}, \quad (9)$$

$$recall = \frac{TP}{TP + FN}, \quad (10)$$

$$F1 = \frac{2 \cdot precision \cdot recall}{precision + recall}, \quad (11)$$

where  $TP$ ,  $FP$ , and  $FN$  represent the number of true positive, false positive and false negative instances respectively. Tables 1 through 3 show the performance evaluation results for layer 2. Specifically, tables 1 and 2 give the scores of text and visual models respectively, and table 3 shows the results of the fused model. As shown in the tables, the fused model outperforms the single-source models. The visual-text model approach achieves a 7% improvement over the text model and a 36% over the visual model. Another observation is that the text model outperforms the visual model. This is because the text information at layer 2 shows a stronger pattern than that of visual information. For example, there is a high possibility that the text files describing images of Cat1 contain the

key "hurricane", while the text files belonging to Cat2 contain the words "oil" and "spill". However, the visual contents of the corresponding images are more abstract and complicated, especially when many categories and subjects are involved. Therefore, a higher weight is assigned to text features at layer 2.

The advantages of text features diminish gradually as the categories are further classified into specific subjects since there is not a strong distinction among those text files in the same category. On the other hand, the visual features demonstrate their superior characteristics for extracting visual patterns when there are fewer subjects involved. Therefore, the weight for visual features increases at layer 3. Tables 4 through 6 contain the subject classification results of layer 3. Specifically, table 4 and table 5 present the scores of text and visual models respectively, and table 6 shows the performance of the combined model. The categorization results of layer 2 enhance the power of visual model at layer 3. The final F1 score of the whole classification framework is 83%, which is 10% and 5% higher than the visual and text models respectively. Although the performance of layer 3 is not as good as layer 2 due to the error propagation problem, the overall experimental results demonstrate the advantages of the data fusion method based on MCA as well as the effectiveness of the hierarchical classification approach.

Categories	Precision	Recall	F1
Cat1	0.98276	0.99415	0.98839
Cat2	0.82698	0.91743	0.86625
Cat3	0.73662	0.97199	0.80158
Average	0.84879	0.96119	0.88541

Table 1: Performance evaluation for text model (Layer-2).

Categories	Precision	Recall	F1
Cat1	0.4485	0.62405	0.51397
Cat2	0.53	0.69419	0.59718
Cat3	0.60715	0.68908	0.6434
Average	0.52855	0.6691	0.58485

Table 2: Performance evaluation for visual model (Layer-2).

## 7 iPad Application Based on HDIC Framework

The proposed HDIC framework has been utilized in an iPad application developed for enhancing disaster situation reports and facilitating decision making processes. The

Categories	Precision	Recall	F1
Cat1	0.98825	0.98533	0.98678
Cat2	0.9578	0.90214	0.9291
Cat3	0.94653	0.93277	0.93925
Average	0.96419	0.94008	0.95171

Table 3: Performance evaluation for visual-text model (Layer-2).

Subjects	Precision	Recall	F1
Sub1	0.87877	0.86807	0.86899
Sub2	0.85186	0.82548	0.83193
Sub3	0.94494	0.85296	0.89631
Sub4	0.92468	0.93569	0.92994
Sub5	0.64674	0.76712	0.67544
Sub6	0.43501	0.58087	0.49665
Average	0.78033	0.80503	0.78321

Table 4: Performance evaluation for text model (Layer-3).

Subjects	Precision	Recall	F1
Sub1	0.64212	0.84219	0.72403
Sub2	0.65548	0.84511	0.73546
Sub3	0.74588	0.80961	0.76942
Sub4	0.75924	0.74552	0.72477
Sub5	0.7355	0.89078	0.79925
Sub6	0.59502	0.75858	0.66468
Average	0.68887	0.8153	0.73627

Table 5: Performance evaluation for visual model (Layer-3).

Subjects	Precision	Recall	F1
Sub1	0.86667	0.88361	0.8707
Sub2	0.88159	0.81884	0.83733
Sub3	0.97143	0.86151	0.91294
Sub4	0.95505	0.91613	0.93427
Sub5	0.71064	0.9133	0.79285
Sub6	0.60397	0.75129	0.66131
Average	0.83156	0.85745	0.8349

Table 6: Performance evaluation for visual-text model (Layer-3).

implementation of the user interface (UI) is based on the officially supported tools for iOS design and coding, i.e., Apple's Xcode 3 and its built-in Interface Builder and iOS Simulator applications. Figure 4 shows the main interface of the system, where users can browse the classified images associated with a specific situation report. Since it is not the

focus of this paper, the details of the implementation are not introduced.



Figure 4: iPad application based on HDIC framework.

## 8 Conclusions and Future Work

In this paper, an hierarchical disaster image classification scheme based on MSDF and MCA is developed for enhancing disaster situation reports with relevant multimedia data and consequently improve the decision making process in disaster situations. The experimental results show the effectiveness of the proposed method. Furthermore, the proposed HDIC framework has been successfully in an iPad application for aiding EM personnel in disaster emergency response. However, there are several aspects of this algorithm to be improved. First, the hierarchical structure and weighing scheme are fixed for a specific scenario, where an adaptive approach is preferable. Second, the visual features are mainly low-level, and more mid-level features are needed to better describe the content of images. Finally, the range of disaster categories and subjects should be extended to serve more general purposes.

## 9 Acknowledgement

This material is based upon work supported by the U.S. Department of Homeland Security under grant Award Number 2010-ST-062-000039 and the U.S. Department of Homeland Security's VACCINE Center under Award Number 2009-ST-061-CI0001.

## References

[1] H. Abdi and D. Valentin. Multiple correspondence analysis. *Encyclopedia of measurement and statistics*, 2007.

[2] A. D. Amo and M. Farmer. Aided image understanding system. *Fuzzy Information Processing Society, NAFIPS*, pages 1–6, 2008.

[3] C. F. Barnes, S. Member, H. Fritz, and J. Yoo. Hurricane disaster assessments with image driven data mining in high resolution satellite imagery. *IEEE Transactions On Geoscience And Remote Sensing Symposium*, pages 1631–1640, 2007.

[4] H. Bayraktar and B. Bayram. Fuzzy logic analysis of flood disaster monitoring and assessment of damage in se anatolia turkey. *Recent Advances in Space Technologies*, pages 13–17, 2009.

[5] S. C. Chen, S. Sista, M. L. Shyu, and R. L. Kashyap. An indexing and searching structure for multimedia database systems. *SPIE Conference on Storage and Retrieval for Media Databases*, pages 262–270, 2000.

[6] S. S. Durbha, R. L. King, V. P. Shah, and N. H. Younan. Image information mining for coastal disaster management. *Proceedings of IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, pages 342–345, 2007.

[7] U. M. Fayyad and K. B. Irani. On the handling of continuous-valued attributes in decision tree generation. *Machine Learning*, 8:87–102, 1992.

[8] S. T. Gandhe, K. T. Talele, and A. G. Keskar. Image mining using wavelet transform. *Knowledge-Based Intelligent Information and Engineering Systems*, pages 797–803, 2007.

[9] C. C. Hsu and Z. Y. Hong. An intelligent typhoon damage prediction system from aerial photographs. *Knowledge-Based Intelligent Information and Engineering Systems*, pages 747–756, 2007.

[10] L. Lin and M. L. Shyu. Weighted association rule mining for video semantic detection. *International Journal of Multimedia Data Engineering and Management*, 1(1):37–54, Jan.-Mar. 2010.

[11] L. Lin, M. L. Shyu, G. Ravitz, and S. C. Chen. Video semantic concept detection via associative classification. *IEEE International Conference on Multimedia and Expo*, pages 418–421, Jul. 2009.

[12] G. Moser and S. B. Serpico. Classification of high resolution images based on mrf fusion and multiscale segmentation. *Proceedings of IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, pages 277–280, 2008.

[13] A. Najab, I. Khan, and F. Ahmad. Principal component analysis based classification of settlements in satellite images. *Proceedings of the 6th International Conference on Frontiers of Information Technology*, 2009.