

Learning Structured Semantic Embeddings for Visual Recognition

Dong Li¹, Hsin-Ying Lee³, Jia-Bin Huang², Shengjin Wang¹, and Ming-Hsuan Yang³

¹Tsinghua University, ²Virginia Tech, ³University of California, Merced

Abstract

Numerous embedding models have been recently explored to incorporate semantic knowledge into visual recognition. Existing methods typically focus on minimizing the distance between the corresponding images and texts in the embedding space but do not explicitly optimize the underlying structure. Our key observation is that modeling the pairwise image-image relationship improves the discrimination ability of the embedding model. In this paper, we propose the structured discriminative and difference constraints to learn visual-semantic embeddings. First, we exploit the discriminative constraints to capture the intra- and inter-class relationships of image embeddings. The discriminative constraints encourage separability for image instances of different classes. Second, we align the difference vector between a pair of image embeddings with that of the corresponding word embeddings. The difference constraints help regularize image embeddings to preserve the semantic relationships among word embeddings. Extensive evaluations demonstrate the effectiveness of the proposed structured embeddings for single-label classification, multi-label classification, and zero-shot recognition.

1. Introduction

Recent visual recognition methods typically train multi-class classifiers using image datasets labeled with a pre-defined set of *discrete* classes [22, 37, 39]. However, such classifiers are not capable of capturing semantic relationships among visual categories since they are trained in the discrete label space. For example, discrete classifiers treat the three classes *cat*, *dog* and *bicycle* as unrelated and distinct categories. As a result, they cannot encode the fact that the two classes *cat* and *dog* are semantically more similar than that between *cat* and *bicycle*. Furthermore, to recognize a new category, the discrete classifiers need to be re-trained on a sufficient amount of training examples of the new class. The lack of semantic information transfer substantially limits the visual recognition methods to scale up to large numbers of classes.

To address these issues, visual-semantic embedding

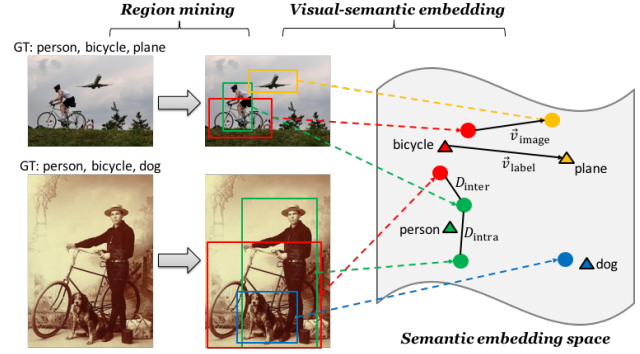


Figure 1. Illustration of the proposed constraints for learning visual-semantic embeddings. Triangles represent label (word) embeddings, and circles represent image embeddings. The visual categories are color-coded. (1) Discriminative constraints (Section 3.2.2) capture the intra- and inter-class relationships of image embeddings (e.g., $D_{intra} < D_{inter}$). (2) Difference constraints (Section 3.2.3) align the difference vector between a pair of image embeddings with that of the corresponding label embeddings (e.g., \vec{v}_{image} and \vec{v}_{label} should be as similar as possible).

models [9, 31, 38] have been proposed to leverage the semantic knowledge from text data. By using a large set of unannotated text data, we can construct a *continuous* and *semantically meaningful* word embedding space [29]. Images can then be mapped into the same semantic space to align the embeddings of their corresponding labels (typically by minimizing ranking losses). The rich semantic relationships from the text data help better recognize visual categories, reduce semantically implausible predictions, and enable zero-shot recognition.

Much effort has been made to learn visual-semantic embeddings by diverse semantic knowledge sources [13, 14, 42, 40]. For example, analogy-preserving embeddings [13] use analogical parallelogram constraints to reflect the relationships between multiple pairs of classes. The analogy relationships can help disambiguate the semantically similar categories. However, learning analogy-preserving embeddings requires manual annotations of attributes and off-the-shelf classifiers to discover a set of analogies. This limits the scalability of the embedding model to handle large numbers of object categories. Structure-preserving con-

straints are also explored to model the neighborhood structure within each modality [42]. Such constraints help improve image-to-text or text-to-image matching by reducing the distance between semantically similar instances. However, they model the neighborhood structure for images and texts separately with two independent regularization terms and thus cannot preserve the semantic relationships between a pair of word embeddings.

In this paper, we propose to learn visual-semantic embeddings by incorporating *discriminative* and *difference* constraints as shown in Figure 1. We exploit the discriminative constraints to explicitly model the intra- and inter-class relationships of image embeddings. Specifically, we explore two types of discriminative constraints (contrastive loss and triplet loss), both of which can help improve the discrimination ability of the embedding model. While discriminative constraints encourage separability for image instances of different categories, there are no constraints on *how* the two instances should be pulled apart. To alleviate these ambiguity issues, we propose the difference constraints to regularize the learning of visual-semantic embedding model. The difference constraints enforce two image embeddings to have similar relative positions with their corresponding label embeddings. Similar to recent work [32, 17, 16], we extend the embedding model to address the *multi-label* scenario where each image may contain multiple labels. Through extensive evaluations, we demonstrate the effectiveness of the proposed structured embeddings for visual recognition. For recognizing seen classes, our method performs better over baseline methods on the CIFAR-10/100 datasets for single-label image classification and achieves competitive performance with the state-of-the-arts on the NUS-WIDE dataset for multi-label image classification. For recognizing unseen classes, our method performs favorably against the state-of-the-arts on the aP&Y and large-scale ImageNet datasets.

We make the following contributions in this work:

First, we exploit the discriminative constraints to learn image-text embeddings using a multi-task learning strategy. The discriminative constraints explicitly model the intra- and inter-class relationships. We show that two types of discriminative constraints can help improve the discrimination ability of the embedding model.

Second, we propose the difference constraints for aligning the difference vectors of image pairs with those of the corresponding label pairs. The difference constraints serve as a regularizer to help learn image embeddings with proper semantic relationships among the various categories.

Third, we present a unified learning formulation that learns visual-semantic embeddings with two additional structured constraints while drawing relations between them. Extensive experimental results show that learning with the two complementary structured constraints signif-

icantly improves visual recognition tasks, including single-label classification, multi-label classification, and zero-shot recognition.

2. Related Work

Visual-semantic embedding. Visual-semantic embedding models relate information from different domains, such as images and texts. Attributes can be used to capture semantic properties shared across different classes [8, 23]. However, attribute-based approaches do not scale up to large amounts of categories due to manually defined attribute ontology and expensive labeling effort. Another line of work leverages neural language models to incorporate semantic knowledge for learning image embeddings [9, 31, 32, 42, 13, 14]. In these approaches, the language model learns semantically meaningful word embeddings from unannotated text data (*e.g.*, [29]). Ranking losses [9, 33] are typically used to learn the image embedding space by constraining the distance between the image embedding and the corresponding word embedding smaller than that between the image embedding and other randomly chosen words. Images are thus projected to nearby positions with their corresponding labels in the semantic space.

Learning embedding with constraints. Other semantic knowledge has also been used to improve embedding models. Examples include analogies [13], taxonomies [14], hierarchies [40] and neighborhood structures [42]. Our work is related to [42] in the aspect of modeling the neighborhood structure of image embeddings. In contrast to learning the transformation layers only, we train our entire network for improved adaptation of visual representations. Moreover, unlike the constraints in [42] that preserve the local neighborhood structure for images and texts *separately*, we regularize that pairs of image embeddings have similar relative positions with their corresponding label embeddings. The proposed difference constraints bear some resemblance with the analogical parallelogram constraints [13], but differ in three aspects. First, we do not rely on costly attribute annotations and off-the-shelf classifiers to discover analogies. Second, the analogy constraints use one pair of classes to help recognize another pair. In contrast, we align the difference vectors between images and labels from the *same* pair of two classes. Third, we incorporate both discriminative and difference constraints into a unified deep learning framework. The contrastive loss [4, 47] or triplet loss [36, 48, 43] has been applied for feature learning. In the context of visual-semantic embedding, we apply either of them as discriminative constraints to improve the embedding baseline using a multi-task learning strategy.

Language grounding methods [18, 19, 24] have been recently proposed for cross-modal tasks (*e.g.*, image-to-text and text-to-image retrieval) by jointly optimizing visual and

semantic embeddings. In this work, we focus on improving the visual model given the pre-trained semantic model.

Convolutional neural networks for visual recognition.

Convolutional neural networks (CNNs) have shown promising results on various visual recognition tasks, *e.g.*, image classification [22, 37, 39]. Recent work addresses the multi-label recognition problem in the discrete label space [11, 41, 25]. For learning visual-semantic embeddings in the general multi-label settings, existing methods either use multiple instance learning [32] or apply off-the-shelf detectors [17, 16] to generate candidate regions for each label. We adopt a different strategy to mine associated regions for each label via a multi-label training procedure. The model pre-trained on the multi-label classification task is also used as an initialization for the subsequent learning steps.

Zero-shot learning. The goal of zero-shot learning is to recognize *unseen* classes without any training. As visual examples of test classes are not available during the training process, auxiliary sources are required to relate the unseen classes with the seen classes. The semantic information sharing across categories can be achieved by attributes [8, 35, 23, 2], word embeddings [9, 31, 38], or a combination of multiple semantic sources [10, 1]. Prior work addresses this problem by learning attribute classifiers [23, 34] or compatibility functions [46, 9, 35, 2, 38]. Recent methods also consider generalized zero-shot learning where test data may come from seen classes and the label space is the union of both seen and unseen classes [3, 45].

3. Approach

Our goal is to learn structured semantic embeddings for visual recognition. We build our learning framework based on deep convolutional neural networks. The pre-trained word embedding model provides continuous vector representations of each image label for training the CNN. To address the multi-label case, we first train the network for multi-label image classification. The learned CNN model is then used to mine top candidate image regions for each label (Section 3.1). Using the mined regions as training instances, we retrain the network to embed image features to the semantic embedding space with our structured constraints (Section 3.2). We include implementation details in Section 3.3.

3.1. Region Mining

Existing visual-semantic embedding methods typically address the single-label setting where each image contains only one semantic label (*e.g.*, ImageNet). This substantially limits the applicability of the learned embedding models for visual recognition as real-world images often contain multiple labels. Similar to recent work [32, 17, 16], we address

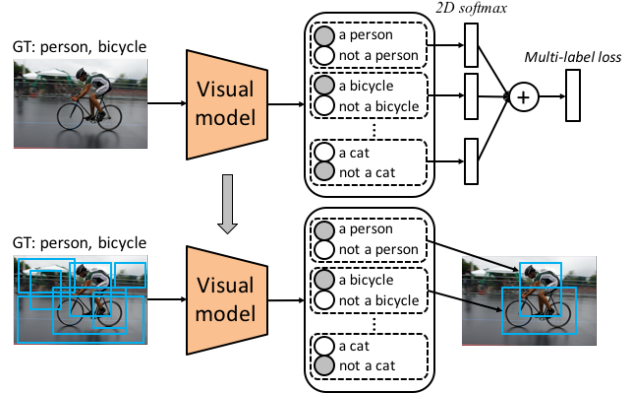


Figure 2. Region mining via multi-label training. First, we use the multi-label loss in [25] to train the network for multi-label image classification. Second, we compute the classification scores of candidate region proposals and select the best-matched region with the highest score for each ground-truth (GT) image label.

this problem by assigning labels to image regions. To this end, we use the multi-label loss [25] to train the network for multi-label image classification. For each ground-truth image label, we compute the classification scores of general region proposals [20] and select the best-matched region with the highest score. Figure 2 illustrates the multi-label classification training process for mining regions that correspond to the image-level labels.

3.2. Structured Semantic Embedding

We use the mined regions as training instances to learn the visual embedding model. We denote the training set as $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where \mathbf{x}_i indicates the i_{th} region instance and y_i indicates the corresponding label. Our goal is to learn a mapping function f_{Θ} that maps from the image space \mathcal{I} to a continuous semantic space \mathcal{S} , $f_{\Theta} : \mathcal{I} \rightarrow \mathcal{S}$, where Θ denotes the network parameters to be optimized. For the word embedding, we exploit the pre-trained word2vec model [29] on the Google News dataset (~ 100 billion words) to generate a 300D vector representation for each label. We denote $s(\cdot)$ as the label embedding function learned by the word2vec model. For the image embedding, we train the CNN to learn Θ by mapping an image to the same 300D space \mathcal{S} . For simplicity, we denote $f(\cdot)$ instead of $f_{\Theta}(\cdot)$ as the image embedding function. Both image and label embeddings are normalized to unit norm.

3.2.1 Baseline Model

Our baseline model aims at projecting image instances to nearby positions with the corresponding labels in the semantic space. We use the ranking loss (1) to learn such an

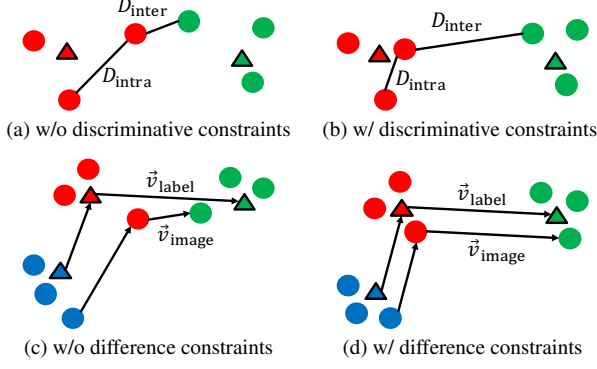


Figure 3. Illustration of the proposed constraints. Triangles represent label embeddings and circles represent image embeddings. The visual categories are color-coded. Discriminative constraints encourage small distance for image instances of the same class and large distance otherwise. Difference constraints align the difference vectors of image embeddings with those of label embeddings.

embedding,

$$L_R(\mathbf{x}_i, y_i) = \sum_{y \neq y_i} \max(0, m + d(f(\mathbf{x}_i), s(y_i)) - d(f(\mathbf{x}_i), s(y))), \quad (1)$$

where we measure the similarity between image and label embeddings based on the cosine distance, *i.e.*, $d(f(\mathbf{x}), s(y)) = 1 - f(\mathbf{x}) \cdot s(y)$.

3.2.2 Discriminative Constraints

The ranking loss (1) optimizes only the distance between the image and label embeddings. However, it does not capture the relationships *among* image embeddings. This may lead to a small margin between images of different classes (see Figure 3(a) for an illustrative example) and limit the discrimination ability of the learned image embeddings. We propose to explicitly model the intra-class and inter-class relationships of image embeddings. Specifically, we apply two alternative discriminative constraints to improve the baseline embedding model. First, the contrastive loss (2) encourages small distance of two images from the same class and large distance otherwise. Second, the triplet loss (3) enforces distance between a reference image and an image from the same class to be smaller than that between the reference image and an image from a different class. With discriminative constraints, image embeddings of the same class are more compact and those of different classes are easier to distinguish, as shown in Figure 3(b). The contrastive loss function is of the form:

$$L_C(\mathbf{x}_i, y_i, \mathbf{x}_j, y_j) = \mathbf{1}_{(y_i=y_j)} d(f(\mathbf{x}_i), f(\mathbf{x}_j)) + \mathbf{1}_{(y_i \neq y_j)} \max(m - d(f(\mathbf{x}_i), f(\mathbf{x}_j)), 0), \quad (2)$$

where $\mathbf{1}(\cdot)$ is the indicator function. The triplet loss function is defined as follows:

$$L_T(\mathbf{x}_i, y_i, \mathbf{x}_j, y_j, \mathbf{x}_k, y_k) = \max(0, m + d(f(\mathbf{x}_i), f(\mathbf{x}_j)) - d(f(\mathbf{x}_i), f(\mathbf{x}_k))), \quad (3)$$

where \mathbf{x}_i denotes the reference image, \mathbf{x}_j is an image from the same class ($y_i = y_j$), and \mathbf{x}_k from a different class ($y_i \neq y_k$).

3.2.3 Difference Constraints

While discriminative constraints enforce image instances of different categories to be distant, there are no constraints on *how* the two instances should be pulled apart. To regularize the learning of visual-semantic embedding model, we propose to align the difference vectors of image pairs with those of label pairs. The difference constraints are capable of preserving the semantic relationships among the label embeddings. Figure 3(c) and (d) illustrate the effect of learning image embeddings with and without using the difference constraints.

We formulate the difference constraints as follows:

$$L_D(\mathbf{x}_i, y_i, \mathbf{x}_j, y_j) = \|(f(\mathbf{x}_i) - f(\mathbf{x}_j)) - (s(y_i) - s(y_j))\|_2^2, \quad (4)$$

where $f(\mathbf{x}_i) - f(\mathbf{x}_j)$ indicates the difference vector of the two image embeddings and $s(y_i) - s(y_j)$ indicates the difference vector for their corresponding label embeddings.

3.2.4 Objective Function

We combine the embedding baseline and two extra structured constraints in a unified learning formulation. The overall training objective function can be represented as:

$$\min_{\Theta} \frac{w}{2} \|\Theta\|^2 + \lambda_1 \sum_i L_R(\mathbf{x}_i, y_i) + \lambda_2 \sum_{i,j} L_C(\mathbf{x}_i, y_i, \mathbf{x}_j, y_j) + \lambda_3 \sum_{i,j} L_D(\mathbf{x}_i, y_i, \mathbf{x}_j, y_j), \quad (5)$$

or

$$\min_{\Theta} \frac{w}{2} \|\Theta\|^2 + \lambda_1 \sum_i L_R(\mathbf{x}_i, y_i) + \lambda_2 \sum_{i,j,k} L_T(\mathbf{x}_i, y_i, \mathbf{x}_j, y_j, \mathbf{x}_k, y_k) + \lambda_3 \sum_{i,j} L_D(\mathbf{x}_i, y_i, \mathbf{x}_j, y_j), \quad (6)$$

where Θ is the parameters of the image embedding function, *i.e.*, the network weights, and $w = 0.0005$ represents the constant weight decay. The weights $\lambda_1, \lambda_2, \lambda_3$ balance these constraints.

3.3. Implementation Details

As shown in Figure 4, we build a two-branch (for contrastive loss) or three-branch (for triplet loss) network to learn image embeddings. We use the AlexNet [22] (for CIFAR-10/100) and GoogLeNet [39] (for the other datasets) pre-trained on the ImageNet 2012 classification task as our base network architectures. Each base network shares the same architecture and parameter weights. We add a linear transformation layer and a normalization layer with randomly initialized parameters on top of the output of each base network. The new transformation layer projects image features to the 300D embeddings. Both image and label embeddings are normalized to unit norm. Since semantic labels may indicate scenes/events/objects in the image, we use the general region proposal method [20] to collect candidate regions. Around 1,000 initial region proposals are generated for each image. We then constrain the region width/height to be at least 0.3 of the image width/height and the aspect ratio to be within the range [0.25, 4].

We use the Caffe toolbox [15] to train CNNs with a Tesla K40 GPU. Since optimizing over all pairs or triplets of instances is computationally infeasible, we randomly sample image instances for training. For the two-branch network, we use equal amounts of image pairs from the same and different classes in a batch. For the three-branch network, we use 20% image instances from the same class of the reference image and the rest from different classes in a batch. We set the initial learning rate to 0.001 with a step decay policy and the momentum to 0.9. We set the margins $m = 0.1$ for the ranking baseline (1) and $m = 1.0$ for both the contrastive (2) and triplet (3) losses in our experiments. The optimal balance weights $\lambda_1, \lambda_2, \lambda_3$ may be different for different datasets. We include the detailed analysis of hyperparameters in the supplementary material. The source code and pre-trained models will be made publicly available.

4. Experimental Results

In this section, we present extensive experimental results to demonstrate the effectiveness of our structured embeddings for visual recognition, including single-label classification, multi-label classification, and zero-shot recognition. We also analyze the contributions of individual components of our approach.

4.1. Datasets

Single-label classification. We first use the CIFAR-10/100 datasets [21] to validate the effectiveness of the proposed constraints. The CIFAR-10 dataset is composed of 10 categories of images with 50,000 training images, and 10,000 testing images. The CIFAR-100 dataset consists of 100 categories. There are 600 images for each category (500 for training and 100 for testing). We evaluate our

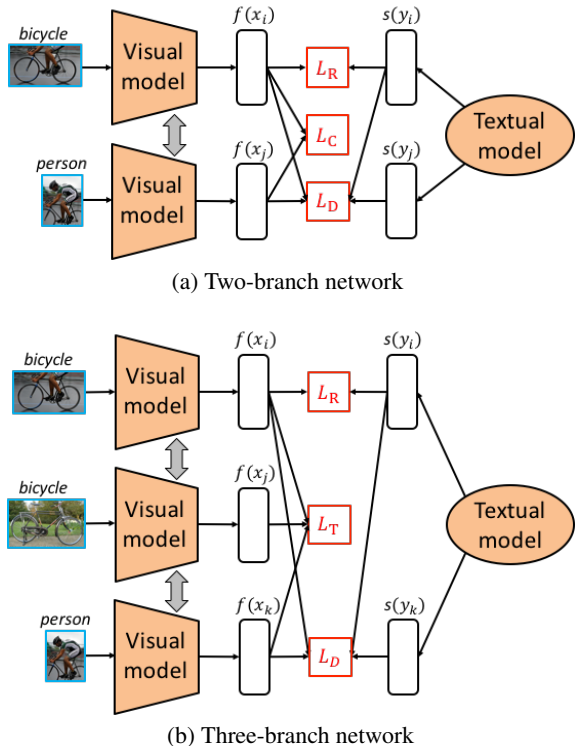


Figure 4. Overview of our network architectures for learning structured embeddings. (a) The two-branch network takes pairs of training instances as input. (b) The three-branch network takes triplets of training instances as input. The base networks share the same architecture and parameter weights. See text for more details of loss functions: L_R (1), L_C (2), L_T (3), L_D (4).

method for image classification on the large-scale ImageNet 2012 dataset with 1,000 labels [6]. We use the default *train* set to train the embedding model with structured constraints and test on the *validation* set.

Multi-label classification. For multi-label image classification, we use the NUS-WIDE dataset [5]. This dataset contains 209,347 images and 81 semantic concepts. We use the train/test split as in [11] and use a subset of 150,000 images for training and the rest of the images for testing.

Zero-shot recognition. For zero-shot recognition, we use the aPascal & aYahoo (aP&Y) [8] and ImageNet 2010 [6] datasets. We follow the standard split on the aP&Y dataset. The Pascal set serves as training data, and the Yahoo set as test data for evaluation. For the ImageNet 2010 dataset, we use the 800/200 split of the 1,000 classes as in [9]: training the embedding model using 800 classes, and inferring image labels using the rest 200 classes.

4.2. Single-Label Classification

Evaluations on CIFAR-10/100. We train a multi-class linear SVM based on the learned 300D embedding features

Table 1. Mean classification accuracy (%) on the CIFAR-10/100 datasets using different constraints. **Red color** and **blue color** indicate the best and second best performing algorithms, respectively.

Rank	Constraints			CIFAR-10	CIFAR-100
	Contrastive	Triplet	Difference		
✓				81.1	62.4
	✓			73.5	42.2
		✓		84.2	53.8
			✓	80.0	58.2
✓	✓			83.1	62.7
✓		✓		86.9	62.5
✓			✓	84.0	61.0
✓	✓		✓	84.6	63.2
✓		✓	✓	87.8	64.3

Table 2. Comparisons of visual-semantic models for image classification on the ImageNet 2012 datasets with 1,000 classes in terms of flat hit@k metrics.

Models	$k = 1$	$k = 2$	$k = 5$	$k = 10$
Norouzi et al. [31]	54.3	61.9	68.0	71.6
Frome et al. [9]	54.9	66.9	78.4	85.0
Rank + Contrastive + Difference	60.7	72.1	82.6	87.9
Rank + Triplet + Difference	57.3	67.7	77.8	83.3

for classification. Table 1 shows the mean classification accuracy using the models trained with different constraints. Combined with discriminative constraints only, we achieve higher accuracy than the baseline model, *e.g.*, a 5.8% gain on CIFAR-10 (86.9% vs. 81.1%). The results demonstrate the effectiveness of discriminative constraints for learning discriminative image embeddings to distinguish visual categories. Combined with difference constraints only, we obtain 2.9% improvement on CIFAR-10 but slightly worse results on CIFAR-100. The intra- and inter-class relationships among different classes are important for image classification. Without explicitly modeling such relationships, difference constraints are not sufficient to distinguish visually similar categories. Combining both constraints, we further improve the classification accuracy on both datasets. The full model trained with triplet and difference losses outperforms the ranking baseline by 6.7% on CIFAR-10 and 1.9% on CIFAR-100. The results validate that our two complementary structured constraints help improve the embedding model for visual recognition. We also use the nearest neighbor classifier where the class label is inferred with the smallest distance between the image and all the candidate label embeddings. We obtain similar performance with SVM (87.0% on CIFAR-10 and 63.8% on CIFAR-100 using the model trained with triplet and different losses).

Evaluations on ImageNet. Our model trained with the contrastive and difference losses achieves 60.7% mean classification accuracy (*i.e.*, flat hit@1), showing a 10.6% relative improvement over the DeVISE method [9]. The results demonstrate the effectiveness of our structured embeddings for large-scale image recognition.

4.3. Multi-Label Classification

We train one-versus-all SVM classifiers for multi-class classification. At test time, we extract one feature vector for each entire image and compute the prediction scores of each class by those SVM classifiers. We also use the setting with region proposals for inference, but find that the predicted regions are not accurate for multi-label classification. Following the same evaluation protocols [11, 32, 41], we generate k (*e.g.*, $k = 3$) highest ranked labels for each test image and then compute the per-class precision (C-P), per-class recall (C-R), per-class F1 (C-F1), overall precision (O-P), overall recall (O-R), and overall F1 (O-F1) scores. The mean average precision (mAP)@N is also used in the recent work [41]. However, we note that these scores are computed based on a *fixed* number of predicted labels for each image. Such metrics are not sufficiently accurate as each image may have a different number of labels. In light of this, we also report the widely used mean average precision (mAP) [7] measure.

Comparisons to the state-of-the-art methods. We compare the proposed visual-semantic embedding approach with the state-of-the-art methods for multi-label image classification, including metric learning [26], multi-edge graph [27], KNN [5], cross-modal ranking [44], WARP [11], MIE [32] and CNN-RNN [41] methods. Table 3 shows quantitative results on the NUS-WIDE dataset. Overall, compared to the methods without learning visual-semantic embeddings [26, 27, 5, 11, 41], our embedding baseline achieves notable improvements in terms of different metrics, *e.g.*, a 10% gain in O-F1 over the KNN baseline [5]. We attribute the performance improvement to the rich semantic relationships among word embeddings, which help better recognize visual categories. Our full model trained with contrastive and difference losses performs favorably against the previous visual-semantic methods [32, 44], *e.g.*, a 2.7% gain over [32] in C-F1 and 37.7% gain over [44] in mAP@10. The results demonstrate the effectiveness of our structured constraints for learning discriminative embedding model. In addition, our full model outperforms the state-of-the-art method [41] by 20% in terms of mAP@10. The results suggest that we get accurate ranked lists in the top 10 predictions for each test image.

Contributions from individual components. We also show the relative contributions of the mined regions and the proposed constraints in Table 3. Our embedding baseline with mined regions achieves 42.4% mAP, significantly outperforming those with random regions or entire images. The performance gain comes from the multi-label training procedure. Compared to the embedding baseline with only the ranking loss, our model learned with both constraints achieves higher performance in terms of all the metrics (*e.g.*, 3.7% improvement in mAP). This demonstrates that

Table 3. Comparisons of image classification performance on the NUS-WIDE dataset. The precision/recall/F1 scores are computed with $k = 3$ predicted labels per image. The ranking loss is enabled in all of our embedding models.

Methods		C-P	C-R	C-F1	O-P	O-R	O-F1	mAP@10	mAP
Li et al. [26]		-	-	-	-	-	21.3	-	-
Liu et al. [27]		-	-	-	35.0	37.0	36.0	-	-
Chua et al. [5]		32.6	19.3	24.3	42.9	53.4	47.6	-	-
Gong et al. [11]		31.7	35.6	33.5	48.6	60.5	53.9	-	-
Wang et al. [41]		40.5	30.4	34.7	49.9	61.7	55.2	56.1	-
Wu et al. [44]		-	-	-	-	-	-	40.3	-
Ren et al. [32]		37.7	40.2	38.9	52.2	65.0	57.9	-	-
W/O region mining	W/ random regions	29.1	3.7	6.6	29.1	35.6	32.0	35.6	3.5
	W/ entire images	33.0	37.6	35.2	50.8	62.0	55.8	73.7	37.0
W/ region mining	Contrastive								
	Triplet								
	Difference								
		35.7	42.9	39.0	52.3	63.9	57.5	76.4	42.4
	✓	35.8	43.3	39.2	52.2	63.8	57.4	76.3	42.5
		38.0	41.2	39.5	52.3	63.8	57.5	76.6	42.2
		36.3	44.0	39.8	52.4	64.0	57.6	76.6	43.2
	✓	38.9	44.6	41.6	52.9	64.6	58.2	77.6	46.1
		38.0	40.9	39.4	52.7	64.3	57.9	77.2	43.9

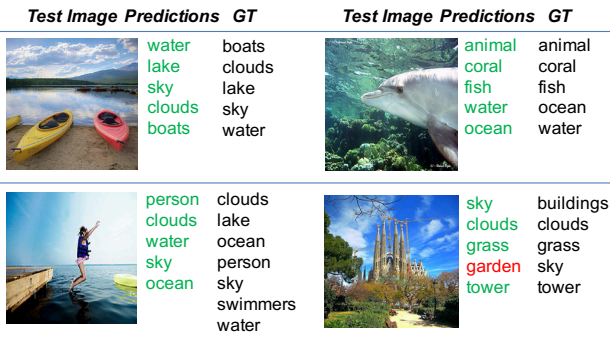


Figure 5. Examples of classification results obtained by our embedding model on the NUS-WIDE dataset. We show top 5 predicted labels for each test image. Green texts indicate correct predictions and red texts indicate wrong predictions. Note that the predictions are ranked and ground-truth labels are not ranked.

our structured constraints help improve the discrimination ability of image embeddings. We observe that there is no significant improvement if we learn the embedding model with discriminative or difference constraints individually. The results show again that the two constraints are complementary for learning visual-semantic embeddings.

Qualitative results. Figure 5 shows examples of prediction results on the test set. Our embedding model makes correct predictions for different semantic concepts, *e.g.*, objects and scenes. The lower right image in Figure 5 shows a typical failure case. Our model incorrectly predicts *garden* that is visually and semantically similar to the ground-truth label of *grass*. While such predictions are incorrect, they are semantically plausible. We refer the readers to the supplementary material for more results.

4.4. Zero-Shot Recognition

One of the critical applications of visual-semantic embedding is zero-shot recognition. For each test image, we first extract the 300D image embedding based on the learned model. We then infer the class label using nearest neighbor search.

Comparisons to the state-of-the-art methods. We compare the proposed approach with the state-of-the-art methods¹ for zero-shot recognition on the aP&Y dataset in Table 4. With word embedding only, we achieve higher accuracy over the existing word embedding based methods, *e.g.*, outperforming DeViSE [9] by 6.7%. The performance gain can be explained by our structured constraints. To combine word and attribute embeddings, we train an additional embedding model with the ranking loss based on the normalized 64D attribute vectors. The word and attribute embeddings are then concatenated for classification using the nearest neighbor search. With the combined embeddings, we obtain competitive results with the state-of-the-art algorithms. Note that our method does not rely on the ground-truth object bounding boxes for training the embedding model. The recent work [50] improves the zero-shot recognition performance by adapting the learned similarity functions using the test data. In contrast, we leverage the nearest neighbor classifier for each test image individually and do not impose additional assumptions on the test data.

Contributions from individual components. We also show the effect of the proposed constraints in Table 4. Our full models consistently outperform the embedding baseline model with the ranking loss by 8.4% with word embedding

¹The results of [31, 9, 1] are from [45] with the ResNet features.

Table 4. Zero-shot image classification accuracy (%) on the aP&Y dataset. The ranking loss is enabled in all of our embedding models.

Semantic sources		Methods	Accuracy
Words	Attributes		
	✓	Romera-Paredes et al. [35]	27.3
	✓	Lampert et al. [23]	38.2
	✓	Zhang et al. [49]	50.4
	✓	Bucher et al. [2]	53.2
✓		Norouzi et al. [31]	25.9
✓		Frome et al. [9]	35.4
✓	✓	Akata et al. [1]	32.0
✓		Contrastive	33.7
✓		Triplet	40.7
✓		Difference	34.3
✓			37.4
✓			42.1
✓			40.2
✓	✓		47.5
✓	✓		48.0
✓	✓		47.8
✓	✓		49.0
✓	✓		54.7
✓	✓		51.1

Table 5. Comparisons of zero-shot recognition on the ImageNet 2010 datasets in terms of flat hit@5 accuracy (%).

Models	200 labels	1,000 labels
Norouzi et al. [31]	28.5	-
Frome et al. [9]	31.8	9.0
Rohrbach et al. [34]	34.8	-
Mensink et al. [28]	35.7	1.9
Fu et al. [10]	41.0	-
Mukherjee et al. [30]	45.7	-
Huang et al. [12]	48.2	-
Rank + Contrastive + Difference	46.2	12.4
Rank + Triplet + Difference	45.0	11.3

only and 7.2% with combined embeddings. Our models with either discriminative or difference constraints also improve the embedding baseline. Combining the two constraints brings further improvement, which shows the complementary nature of the two structured constraints.

Error analysis. We analyze per-class recognition results of our method with the confusion matrix in Figure 6(a). The majority of errors comes from confusion with semantically similar categories, *e.g.*, *zebra* and *donkey*. This is because words with similar semantics are embedded at close positions in the space, as shown in Figure 6(b).

Evaluations on ImageNet. We show the capacity of our visual-semantic embedding models for zero-shot recognition when there are large amounts of class labels. Table 5 shows performance comparisons in terms of flat hit@5 accuracy. Our models learned with the proposed structured

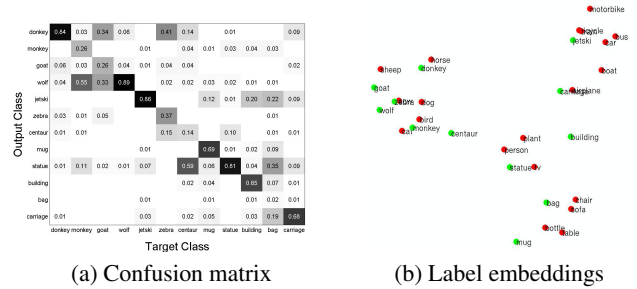


Figure 6. Error analysis of zero-shot recognition by our method on the aP&Y dataset. (a) Confusion matrix computed based on the per-class prediction results using our full model. Most errors come from semantically similar classes. (b) The label embeddings of both seen classes from Pascal (red) and unseen classes from Yahoo (green). The object classes with similar semantics are close in the embedding space.

constraints achieve comparative performance with the state-of-the-art methods. When classifying test data into the joint label space of both seen and unseen classes, we also achieve better accuracy (a 3.4% gain over DeViSE [9]). This indicates that our models have less bias toward training classes than the previous methods.

5. Conclusions

In this paper, we propose to incorporate two structured constraints for learning visual-semantic embeddings. Discriminative constraints model the intra- and inter-class relationships and difference constraints serve as a regularizer to preserve the semantic relationships among word embeddings. Quantitative results show that the two constraints are complementary and crucial for improving visual recognition. Our method is simple, flexible, and easily applicable to large amounts of categories since we do not rely on costly bounding box annotations. Experimental evaluations on multiple datasets including the large-scale ImageNet dataset demonstrate the effectiveness of our embedding model with structured constraints for image classification and zero-shot recognition. In the future work, we plan to jointly learn the visual and textual embeddings and explore additional applications, *e.g.*, object localization using the visual-semantic embedding model.

References

- [1] Z. Akata, S. Reed, D. Walter, H. Lee, and B. Schiele. Evaluation of output embeddings for fine-grained image classification. In *CVPR*, 2015. 3, 7, 8
- [2] M. Bucher, S. Herbin, and F. Jurie. Improving semantic embedding consistency by metric learning for zero-shot classification. In *ECCV*, 2016. 3, 8
- [3] W.-L. Chao, S. Changpinyo, B. Gong, and F. Sha. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *ECCV*, 2016. 3
- [4] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *CVPR*, 2005. 2
- [5] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng. Nus-wide: a real-world web image database from national university of singapore. In *ACM CIVR*, 2009. 5, 6, 7
- [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 5
- [7] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2):303–338, 2010. 6
- [8] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth. Describing objects by their attributes. In *CVPR*, 2009. 2, 3, 5
- [9] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al. Devise: A deep visual-semantic embedding model. In *NIPS*, 2013. 1, 2, 3, 5, 6, 7, 8
- [10] Z. Fu, T. Xiang, E. Kodirov, and S. Gong. Zero-shot object recognition by semantic manifold distance. In *CVPR*, 2015. 3, 8
- [11] Y. Gong, Y. Jia, T. Leung, A. Toshev, and S. Ioffe. Deep convolutional ranking for multilabel image annotation. In *ICLR*, 2014. 3, 5, 6, 7
- [12] C. Huang, C. C. Loy, and X. Tang. Local similarity-aware deep feature embedding. In *NIPS*, 2016. 8
- [13] S. J. Hwang, K. Grauman, and F. Sha. Analogy-preserving semantic embedding for visual object categorization. In *ICML*, 2013. 1, 2
- [14] S. J. Hwang and L. Sigal. A unified semantic embedding: Relating taxonomies and attributes. In *NIPS*, 2014. 1, 2
- [15] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *ACM MM*, 2014. 5
- [16] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015. 2, 3
- [17] A. Karpathy, A. Joulin, and F. F. Li. Deep fragment embeddings for bidirectional image sentence mapping. In *NIPS*, 2014. 2, 3
- [18] R. Kiros, R. Salakhutdinov, and R. S. Zemel. Multimodal neural language models. In *ICML*, 2014. 2
- [19] S. Kottur, R. Vedantam, J. M. F. Moura, and D. Parikh. Visual word2vec (vis-w2v): Learning visually grounded word embeddings using abstract scenes. In *CVPR*, 2016. 2
- [20] P. Krähenbühl and V. Koltun. Geodesic object proposals. In *ECCV*, 2014. 3, 5
- [21] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Master’s thesis, Computer Science Department, University of Toronto, 2009. 5
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1, 3, 5
- [23] C. H. Lampert, H. Nickisch, and S. Harmeling. Attribute-based classification for zero-shot visual object categorization. *TPAMI*, 36(3):453–465, 2014. 2, 3, 8
- [24] A. Lazaridou, N. T. Pham, and M. Baroni. Combining language and vision with a multimodal skip-gram model. In *NAACL*, 2015. 2
- [25] D. Li, J.-B. Huang, Y. Li, S. Wang, and M.-H. Yang. Weakly supervised object localization with progressive domain adaptation. In *CVPR*, 2016. 3
- [26] J. Li, X. Lin, X. Rui, Y. Rui, and D. Tao. A distributed approach toward discriminative distance metric learning. *TNNLS*, 26(9):2111–2122, 2015. 6, 7
- [27] D. Liu, S. Yan, Y. Rui, and H.-J. Zhang. Unified tag analysis with multi-edge graph. In *ACM MM*, 2010. 6, 7
- [28] T. Mensink, J. Verbeek, F. Perronnin, and G. Csurka. Metric learning for large scale image classification: Generalizing to new classes at near-zero cost. In *ECCV*, 2012. 8
- [29] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013. 1, 2, 3
- [30] T. Mukherjee and T. Hospedales. Gaussian visual-linguistic embedding for zero-shot recognition. In *EMNLP*, 2016. 8
- [31] M. Norouzi, T. Mikolov, S. Bengio, Y. Singer, J. Shlens, A. Frome, G. S. Corrado, and J. Dean. Zero-shot learning by convex combination of semantic embeddings. In *ICLR*, 2014. 1, 2, 3, 6, 7, 8
- [32] Z. Ren, H. Jin, Z. Lin, C. Fang, and A. Yuille. Multi-instance visual-semantic embedding. *arXiv preprint arXiv:1512.06963*, 2015. 2, 3, 6, 7
- [33] Z. Ren, H. Jin, Z. Lin, C. Fang, and A. Yuille. Joint image-text representation by gaussian visual-semantic embedding. In *ACM MM*, 2016. 2
- [34] M. Rohrbach, M. Stark, and B. Schiele. Evaluating knowledge transfer and zero-shot learning in a large-scale setting. In *CVPR*, 2011. 3, 8
- [35] B. Romera-Paredes and P. Torr. An embarrassingly simple approach to zero-shot learning. In *ICML*, 2015. 3, 8
- [36] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015. 2
- [37] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 1, 3
- [38] R. Socher, M. Ganjoo, C. D. Manning, and A. Ng. Zero-shot learning through cross-modal transfer. In *NIPS*, 2013. 1, 3
- [39] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 1, 3, 5
- [40] I. Vendrov, R. Kiros, S. Fidler, and R. Urtasun. Order-embeddings of images and language. In *ICLR*, 2016. 1, 2
- [41] J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, and W. Xu. Cnn-rnn: A unified framework for multi-label image classification. In *CVPR*, 2016. 3, 6, 7
- [42] L. Wang, Y. Li, and S. Lazebnik. Learning deep structure-preserving image-text embeddings. In *CVPR*, 2016. 1, 2
- [43] X. Wang and A. Gupta. Unsupervised learning of visual representations using videos. In *ECCV*, 2015. 2
- [44] F. Wu, X. Jiang, X. Li, S. Tang, W. Lu, Z. Zhang, and Y. Zhuang. Cross-modal learning to rank via latent joint representation. *TIP*, 24(5):1497–1509, 2015. 6, 7
- [45] Y. Xian, Z. Akata, and B. Schiele. Zero-shot learning – the Good, the Bad and the Ugly. In *CVPR*, 2017. 3, 7
- [46] Y. Xian, Z. Akata, G. Sharma, Q. Nguyen, M. Hein, and B. Schiele. Latent embeddings for zero-shot classification. In *CVPR*, 2016. 3
- [47] S. Zhang, Y. Gong, J.-B. Huang, J. Lim, J. Wang, N. Ahuja, and M.-H. Yang. Tracking persons-of-interest via adaptive discriminative features. In *ECCV*, 2016. 2
- [48] X. Zhang, F. Zhou, Y. Lin, and S. Zhang. Embedding label structures for fine-grained feature representation. In *CVPR*, 2016. 2
- [49] Z. Zhang and V. Saligrama. Zero-shot learning via joint latent similarity embedding. In *CVPR*, 2016. 8
- [50] Z. Zhang and V. Saligrama. Zero-shot recognition via structured prediction. In *ECCV*, 2016. 7