

Semi-supervised Zero-Shot Learning by a Clustering-based Approach

Seyed Mohsen Shojaee and Mahdieh Soleymani Baghshah

Sharif University of Technology
Tehran, Iran

mshojaee@ce.sharif.edu, soleymani@sharif.edu

Abstract

In some of object recognition problems, labeled data may not be available for all categories. Zero-shot learning utilizes auxiliary information (also called signatures) describing each category in order to find a classifier that can recognize samples from categories with no labeled instance. In this paper, we propose a novel semi-supervised zero-shot learning method that works on an embedding space corresponding to abstract deep visual features. We seek a linear transformation on signatures to map them onto the visual features, such that the mapped signatures of the seen classes are close to labeled samples of the corresponding classes and unlabeled data are also close to the mapped signatures of one of the unseen classes. We use the idea that the rich deep visual features provide a representation space in which samples of each class are usually condensed in a cluster. The effectiveness of the proposed method is demonstrated through extensive experiments on four public benchmarks improving the state-of-the-art prediction accuracy on three of them.

Introduction

Zero-shot learning (Larochelle, Erhan, and Bengio 2008; Palatucci et al. 2009; Lampert, Nickisch, and Harmeling 2009; Farhadi et al. 2009) is an extension to the conventional supervised learning scenario that releases the assumption that ample labeled instances is available for all categories. It addresses the recognition problem when no labeled instance is available for some classes. Instead, some sort of description that is called *class signature* is available for all categories. Example of class signatures include human-annotated discriminative attributes or textual description of the categories. The problem addressed by zero-shot learning rises naturally in practice wherever it is not feasible to acquire abundant labeled instances for all the categories (e.g., fine-grained classification problems). To describe the task more precisely, in the training phase, labeled instances for some categories which are called seen classes are provided while for other categories, called unseen ones, there is no labeled instance available. In the test phase, unlabeled instances should be classified into seen or unseen categories. In this work, we focus on the most popular version of zero-

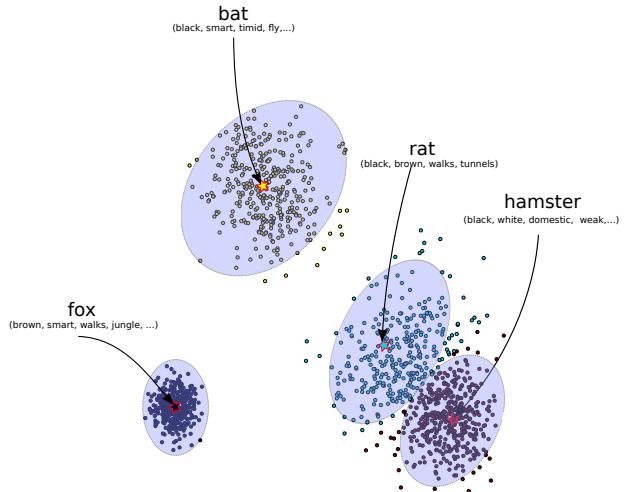


Figure 1: In our model we map class signatures close to all instances of the class for seen classes and close to instances in a cluster for unseen ones. Three different clusters are shown here by circles with different colors and class signatures are mapped to center of clusters shown by stars of the same color

shot recognition in which test instances belong only to unseen categories.

Most existing methods for zero-shot learning focus on using labeled images to learn a compatibility function indicating how similar an image is to each label embedding (Akata et al. 2015b; Romera-Paredes and Torr 2015; Zhang and Saligrama 2015b). Each instance will then be labeled with the category having the most compatible signature. On the other hand, recent advances in deep neural networks provide rich visual features with high discrimination capability (Simonyan and Zisserman 2014). We will show in Section through experiments that the space of deep visual features is a rich space in which instances of different categories usually form natural clusters. Little attention has been paid to exploiting this property of visual features in the context of zero-shot learning.

In this paper, we propose Joint Embedding and Clustering (JEaC), a semi-supervised zero-shot learning method in which both supervised information from labeled samples

and unsupervised information from unlabeled samples are utilized and the assignment of unlabeled samples to unseen classes is jointly learned with a linear embedding from class signatures to the space deep image features (Figure 1). The linear mapping is learned so that the mapped signature of each seen class tends to be representative for samples of the corresponding class and simultaneously assignments of unlabeled samples to unseen classes are leaned such that the mapped signature of each unseen class also tends to be representative of assigned samples to that class. Using unlabeled samples from unseen classes, we can substantially mitigate the *domain shift problem* introduced in (Fu et al. 2014) that impairs the zero-shot recognition performance. We also propose *Independent Embedding and Clustering* (IEAC), a simpler method in which label assignment and class embedding in image space are not learned jointly. Instead, after finding the mapping using only labeled data, a clustering algorithm is used to assign labels to instances of unseen classes.

In Section , we present experimental results on four popular zero-shot classification benchmarks and see that the proposed method outperforms the state-of-the-art methods on these datasets. The rest of paper is organized as follow. In Section , we briefly introduce existing methods for zero-shot learning. In Section , we present our semi-supervised zero-shot learning method. In Section , we report our experiments and finally in Section we conclude.

Related Work

A notable body of work in zero-shot recognition belongs to attribute prediction from images (Lampert, Nickisch, and Harmeling 2009; Yu and Aloimonos 2010; Mahajan, Sellamanickam, and Nair 2011; Wang and Ji 2013; Suzuki et al. 2014). In these methods, the semantic label embeddings are considered to be externally provided attributes. Thus, label embeddings are already available and the task is just to map images to the semantic space, i.e., predicting attributes for the images. Early methods, like (Lampert, Nickisch, and Harmeling 2009), assume independence between attributes and train binary attribute classifiers. Probabilistic graphical models have been utilized to model and/or learn correlations among different attributes (Yu and Aloimonos 2010; Wang and Ji 2013) to improve the prediction. In (Jayaraman and Grauman 2014), a random forest approach has been employed that accounts for unreliability in attribute predictions for final class assignment. In (Akata et al. 2015a), a max-margin objective function is defined for attribute-based image classification.

More recent works exploit bilinear models (Yu et al. 2013; Frome et al. 2013; Norouzi et al. 2014; Zhang and Saligrama 2015b; Romera-Paredes and Torr 2015; Schuurmans and Tg 2015). Several objective functions have been proposed for learning such bilinear models. In (Romera-Paredes and Torr 2015), the sum of the squared error on the label prediction is used and clever regularization terms that compensate undesirable characteristics of this cost function are also utilized. In (Li and Guo 2015; Schuurmans and Tg 2015) a max margin objective function for learning the bilinear mapping is

used. These two methods learn labels for test instances simultaneously and so they differ from almost all of other existing methods in this way. This provides the possibility of leveraging unsupervised information available in test images, for instance as done in (Schuurmans and Tg 2015), by using a Laplacian regularization term that penalizes similar objects assigned to different classes.

Designing label embeddings in multi-class classification is another line of research that can also be used in zero-shot recognition. In (Yu et al. 2013), an objective function is proposed to derive such label embeddings based on information about similarities among categories. A relatively popular embedding for labels is to describe unseen categories as how similar they are to the seen ones. In (Norouzi et al. 2014), the outputs from the softmax layer of a CNN trained on seen categories are used to score similarity between test instances and seen classes. Using these outputs as weights, images are represented in the semantic space as a convex combination of seen class label embeddings. In (Zhang and Saligrama 2015b), a histogram showing seen class proportions is used for label embedding and then a max margin framework is defined to embed images in this space. this work is further extended in (Zhang and Saligrama 2015a) where a supervised dictionary learning formulation is presented to jointly learn embedding for images and label. The idea of combining already available classifiers to create new ones for unseen categories is also used in (Changpinyo et al. 2016) but rather than using seen categories as basis, they define a set of (possibly smaller) *phantom* classes and learn base classifiers on them.

Although most of the works for zero-shot recognition consider attributes as auxiliary information, textual descriptions or name of the categories is also used as class signatures. (Frome et al. 2013) introduce a bilinear model to find the compatibility score of deep visual features and Word2vec (?) representation of class names. (Ba et al. 2015) proposes nonlinear mappings modeled by neural networks on the image and the text inputs to find their compatibility. (Elhoseiny, Saleh, and Elgammal 2013) presents an objective function to predict classifier parameters from textual descriptions. In (Akata et al. 2015b), different label embeddings and also a combination of these different embeddings have been considered along with a bilinear compatibility function. In (Xian et al. 2016), this work is extended further to model nonlinear compatibility functions that can be expressed as a mixture of bilinear models. In (Fu and Sigal 2016), a set of vocabulary much larger than just seen and unseen class names is used and mapping from images to word embeddings is learned by maximizing the margin with respect to all words in the vocabulary; this framework can be used in zero-shot and also supervised and open set learning problems. In (Akata et al. 2016), authors propose to use multiple auxiliary information and also part annotation in image domain to compensate for weaker supervision in textual data. Convolutional and recurrent neural networks have also been used for text embedding in (Scott Reed, Zeynep Akata, Honglak Lee 2016).

The most related methods to our method are (Li and Guo 2015; Kadirov et al. 2015; Schuurmans and Tg 2015) that

are semi-supervised zero-shot learning methods. Here, we briefly specify the differences between these methods and ours. First, we use abstract visual features obtained by deep learning as the semantic space as opposed to (Li and Guo 2015; Schuurmans and Tg 2015). This two methods learn a max margin classifier on the image space classifying both seen and unseen instances while we use a ridge regression to map signatures to the semantic visual space resulting in a much simpler optimization problem to solve. We also explicitly account for domain shift problem in our objective function and thus achieve better results compared to these methods. There are major differences between our work and (Kodirov et al. 2015) that use a dictionary learning scheme in which coding coefficients are considered to be label embeddings (attribute vectors) and a sparse coding objective is used to map images into this representation space. Most importantly, in our method labels of unseen instances are jointly learned with the mapping of the signatures to the semantic space in our objective function while in (Kodirov et al. 2015) the label prediction is accomplished using the nearest neighbor or the label propagation on embeddings of images. Also, we do not need to learn embedding of test instances in the semantic space as opposed to (Kodirov et al. 2015), alternatively we learn just the representation of class signatures in the visual domain.

Proposed Method

In this section, we introduce two zero-shot learning method that use the deep visual features as the semantic space and learn a mapping from class signatures to this semantic space and also predict labels for instances belonging to unseen classes. First, we propose Independent Embedding and Clustering (IEaC), a simple and efficient semi-supervised zero-shot learning method in Section . Then, we further extend our method to jointly learn the embedding of signatures and class assignments of unlabeled samples. We call this method Joint Embedding and Clustering (JEaC). We formulate JEaC as an optimization problem and introduce an iterative method to solve this it in Section . We use IEaC to find a starting point for this optimization procedure (i.e., as an initial labellings for instances of unseen classes).

Notation

Let X , \mathbf{x} , and x denote matrices, column vectors, and scalars respectively. $\|X\|_F^2$ shows the squared Frobenius norm of a matrix and $X_{(i)}$ denotes its i th column. $\mathbf{1}_k$ denotes a column vector whose $k - th$ element is one and is zero everywhere else. Suppose there are n_s seen categories and n_u unseen categories. For each category y , auxiliary information $a_y \in \mathbb{R}^r$ is available. We assume that labels $\{1, \dots, n_s\}$ correspond to seen categories.

Let $X_s \in \mathbb{R}^{d \times N_s}$ and $X_u \in \mathbb{R}^{d \times N_u}$ denote matrices whose columns are seen and unseen images respectively where d is the dimension of image features. $S_s = [a_1, \dots, a_{n_s}]$ presents the matrix of signatures for seen classes. S_u is also defined similarly for unseen classes. $Z_s = [\mathbf{z}_1, \dots, \mathbf{z}_{N_s}]$ contains labels of training data in one-hot encoding format.

Independent Embedding and Clustering

Our first method can be roughly summarized in three steps:

1. Using data from seen classes, we learn a linear mapping from attribute vectors to the semantic space.
2. We find a data clustering using our proposed semi-supervised clustering algorithm.
3. For instances of each cluster, we find the label whose mapped signature in the semantic visual space is the nearest one to the center of that cluster and assign that label to all of these instances.

We use a simple ridge regression to map class signatures to visual features. We intend to find a mapping from class signatures to the deep visual representation space such that each mapped (seen) class signature is close to the samples of that class in this space in average. The linear mapping is found using the following optimization problem:

$$D = \arg \min_D \|X_s - DY_s\|_F^2 + \gamma \|D\|_F^2, \quad (1)$$

where columns of $Y_s \in \mathbb{R}^{r \times n_s}$ are the class signatures of the samples lied in the columns of X_s . This optimization problem is known to have the following closed form solution:

$$D = X_s Y_s^T (Y_s Y_s^T + \gamma I)^{-1}. \quad (2)$$

The parameter γ is determined through cross validation as we will describe precisely in Section .

Here, we intend to find labels for instances belonging to unseen classes. To this end, we want to find a clustering of instances in the space of deep visual features and assign a label to each cluster according to the distance between the center of that cluster and the mapped signature of the unseen classes (i.e., consider the label whose mapped signature is the closest one to the cluster center as the assigned label to the instances of this cluster). To find a better clustering of instances belonging to unseen classes, we can also incorporate labeled instances of seen classes too. The clustering problem over unseen instances, we encountered here, is different from the conventional semi-supervised learning problem (Chapelle, Schölkopf, and Zien 2006). In fact, all labeled data are from seen classes and there is no labeled sample for unseen classes that is due to the special characteristic of zero-shot learning problem. Therefore, here, we propose a semi-supervised learning method which can be seen as an extension to k-means suitable for this problem. We try to find a clustering such that labeled instances tend to be assigned to the corresponding classes and all instances tend to be close to the center of the clusters to which they are assigned:

$$\min_{R, \mu_1, \dots, \mu_k} \sum_{n,k} r_{nk} \|\mathbf{x}_n - \mu_k\|_2^2 + \beta \sum_{n=1}^{N_s} \mathbf{1}(\mathbf{r}_n \neq \mathbf{z}_n), \quad (3)$$

where μ'_k s are cluster centers and $R = [\mathbf{r}_1, \dots, \mathbf{r}_{N_s+N_u}]$ is cluster assignments in one-hot encoding format. The objective function is similar to that of the k-means clustering algorithm but for each labeled instance there is a penalty of β if its assigned cluster number that is different from its label.

Thus, this objective function encourages the first n_s clusters be corresponding to the seen classes.

Parameters β and k can be determined via the cross validation. However, in our experiments, we found out the model is not very sensitive to these so we fix $\beta = 1$ when data is normalized such that $\|x_i\|_1 = 1$. We set $k = (n_s + 2n_u)$, this will allow for two clusters per class for unseen categories. This, to some extent, copes with diversity in instances of a class.

Finally, to assign labels to test instances, we use the mapping D from Eq.(2) to map class signatures to visual features, creating a set of *class representatives* in the visual feature space. We then assign to all instances of a cluster the class label whose representative is the nearest to center of that cluster.

A key distinction between the clustering-based method presented here and other existing methods lies in the nature of the compatibility function. The compatibility function in other works is a similarity measure between each instance and class description that is found independently for different instances. Here, the compatibility function relies strongly on the distribution of instances in the semantic space and the compatibility of a label for an instance is found according to the similarity of the cluster center to which this instance is assigned and the mapped signature of that label. Therefore, by considering the distribution of data points (via clustering) in designing the compatibility function we can reach a more reliable measure. This compatibility function can be plugged in every other method in this way that after final predictions are made by the method, a clustering algorithm is ran on data and then we assign an identical label to all cluster members by majority voting on those predictions.

Although the above method outperforms the state-of-the-art methods on most zero-shot recognition benchmarks, it uses only the instances of the seen classes to find the linear transformation from class signatures to the visual feature space and thus the proposed method may suffer from the domain shift problem introduced in (Fu et al. 2014). To overcome the domain shift problem more substantially, we propose an optimization problem for finding the linear transform from class signatures to the visual feature space that uses instances of both seen and unseen classes.

Joint Embedding and Clustering

In this section, we propose an optimization problem for learning a linear transformation from class signatures to the visual features space such that the mapped signatures are good representatives of the corresponding instances. We intend to learn a transformation such that for the seen classes, the sum of the squared distances of instances from the mapped signature of the corresponding class is minimized. Moreover, for instances of unseen classes, we can find class assignments such that the sum of the squared distances of unseen instances from the mapped signature of classes to which they are assigned is also minimized. The objective

Algorithm 1: Training Procedure of IEaC

```

input :  $X_s, Y_s, Z_s, X_u, S_u$ 
output:  $Z_u$  (label predictions for  $X_u$ )
 $k \in \{1, 2, \dots, n_s + n_u\}$ 
 $n \in \{1, 2, \dots, N_s + N_u\}$ 
Initialize  $\mu_k$  by Eq. (5),  $k = 1, \dots, n_s$ ;
Initialize  $\mu_k$  by kmeans++,  $k = n_s + 1, \dots, n_s + n_u$ ;
repeat
     $c_n \leftarrow \arg \min_i \|x_n - \mu_i\|_2$ ; //cluster assignments
     $\mu_k \leftarrow \sum_n \mathbf{x}_n \mathbb{1}(c_n = k) / \sum_n (\mathbb{1}(c_n = k))$  ;
until convergence to local minimum;
 $D \leftarrow X_s Y_s^T (Y_s Y_s^T + \gamma I)^{-1}$ ;
// array  $l$  maps cluster numbers to labels
 $l[k] \leftarrow \arg \min_j \|\mu_k - (DS_u)_{(j)}\|_2$ ;
 $(Z_u)_{(n)} \leftarrow \mathbf{1}_{l[c_n]}$ ;

```

function of JEaC is formulated as follows:

$$\begin{aligned} \min_{R, D} & \|X_s - DY_s\|_F^2 + \lambda \|X_u - DS_u R^T\|_F^2 + \gamma \|D\|_F^2 \\ \text{s. t. } & R \in \{0, 1\}^{N_u \times n_u}. \end{aligned} \quad (4)$$

The first term in the above optimization problem is identical to Eq.(1) and the second one incorporates unlabeled data for learning the mapping D . By enforcing the signatures to be mapped close to test instances, this term confronts the domain shift problem. In fact, we seek a class assignment for instances of unseen classes such that we can learn a linear transformation on class signature to use the mapped signature of both seen and unseen classes as good representatives for the corresponding instances. The second term can be essentially considered as a clustering objective with two advantages. First, the number of clusters is no longer a parameter and it is determined by the number of unseen classes. Second, the cluster centers are set to be the mapped signatures of test classes.

Optimization

Training Algorithm for IEaC: Optimization of the objective function in Eq. (3) is done by alternating between μ'_i s and R . μ'_i s are updated using:

$$\mu'_i = \frac{\sum_{n=1}^{N_s+N_u} \mathbb{1}(r_{ni} = 1) \mathbf{x}_n}{\sum_{n=1}^{N_s+N_u} \mathbb{1}(r_{ni} = 1)}, \quad (5)$$

R is updated by assigning each instance to the cluster that minimizes the corresponding term in Eq.(3). To initialize μ'_i s, for clusters corresponding to seen classes the centers are set as mean of instances from that class. Centers of other clusters are initialized using k-means++ (Arthur and Vassilvitskii 2007) on unlabeled instances. The overall training algorithm IEaC is presented in algorithm 1

Training algorithm for JEaC: The Eq. (4) is not convex and considering that R is a partitioning of instances, the global optimization requires an exhaustive search over all possible labeling of test data with n_u labels. Therefore, we use a simple coordinate descent method (like k-means). We

Algorithm 2: Training Procedure for JEaC

input : X_s, Y_s, Z_s, X_u, S_u
output: Z_u (label predictions for X_u)
Initialize R by output of Algorithm 1 ;
repeat
 | update D by Eq. (6) ;
 | updata R by Eq. (7) ;
until no element of R changes;
output $Z_u \leftarrow R$

alternate between optimizing R and D while fixing the other. Having fixed the labeling R , the problem becomes a simple multi-task ridge regression which has the following closed-form solution:

$$D = (X_s Y_s^T + \beta X_u R S_u^T)(Y_s Y_s^T + \beta S_u R^T R S_u^T + \gamma I)^{-1}. \quad (6)$$

By fixing D , the optimal R can be achieved via assigning each instance to the closest class representative:

$$r_{ij} = \mathbb{1}[j = \arg \min_k \|X_{u(i)} - D S_{u(k)}\|_2]. \quad (7)$$

Whenever a row of R contains no 1's, i.e. an empty cluster is encountered we assign 2% of instances randomly to that cluster. We continue alternating between updates of D and R till R remains constants, i.e., no label changes. In our experiments, this always happens in less than 20 iterations.

To evade poor local minima, we use a good starting point that initializing R with prediction made by IEaC.

Experiments

In this section, we conduct experiments on the popular benchmarks to obtain results of the proposed method on these benchmarks and compare them with those of the other methods.

Datasets. We evaluate our proposed methods on four popular public benchmarks for zero-shot classification. (1) Animal with Attributes (AwA) (Lampert, Nickisch, and Harmeling 2009). There are images of 50 mammal species in this data set. Each class is described by a single 85-dimensional attribute vector. We use the continuous attributes rather than the binary version as it has proved to be more discriminative in previous works like (Akata et al. 2015b). The train/test split provided by the dataset is used accordingly. (2) aPascal/aYahoo (Farhadi et al. 2009). The 20 categories from Pascal VOC 2008 (Hoiem, Divvala, and Hays 2008) are considered as seen classes and categories from aYahoo are considered to be unseen. As this dataset provides instance level attribute vectors, for class signatures we use the average of the provided instance attributes. (3) SUN Attribute (Patterson et al. 2014). The dataset consists of 717 categories and all images are annotated with 102 attributes, we just use the average attributes among all instances of each categories for our experiments. We use the same train/test spilt as in (Jayaraman and Grauman 2014) where 10 classes have been considered unseen. (4) Caltech UCSD Birds-2011 (CUB) (Wah et al. 2011). This a dataset for fine-grained classification task. There are 200 species of birds where each image

has been annotated with 312 binary attributes. Again, we average over instances to get continuous class signatures. We use the same train/test split as in (Akata et al. 2013) (and many other following works) to make comparison possible.

As our method relies on meaningful structure in visual features domain, we use features from a deep CNN known that are more discriminative than *shallow* features like SIFT or HOG. We report results using 4096-dimensional features from the first fully connected layer of 19 layer VGG network (Simonyan and Zisserman 2014) pre-trained on image-net, provided publicly by (Zhang and Saligrama 2015b).

Testing Cluster Assumption: First, to give evidence for our key assumption of our method that instances from each class usually form a cluster in visual feature domains and to demonstrate effectiveness of our proposed clustering algorithm we design an experiment in which instances from unseen categories are clustered using our proposed clustering algorithm and also the k-means algorithm. Then, each cluster is assigned with a class label based on majority voting on ground truth labels. The number of clusters is set to the number of classes as a natural choice.¹ For the k-means algorithm, we use the implementation available in Scikit-learn library (Pedregosa et al. 2011) and run it with 20 different initializations and report results of that one with the best score. Accuracy of this labeling scheme that is based on clustering is reported in Table 1. These results shows the effectiveness of our proposed clustering method and that the cluster structure assumption in the visual semantic space is usually right.

Cross Validation: To adjust parameters γ and β in Eq. 6 and parameter γ in Eq. 2, we split training data into train and validation sets. We choose a number of categories randomly from training data as validation categories. For each data set, the size of the validation set has the same ratio to the train set as the size of the test categories to the total of the train and the validation one. In our experiments, we used 10-fold cross validation, i.e., average results from ten different validation splits are used to decide on optimal parameters. Once optimal γ and β are determined through the grid search by testing on validation set, the model is then trained on all seen categories.

We summarize our experimental results in Table 2. *IEaC* (*Proposed Clustering*) corresponds to the method presented in Section . *IEaC* (*K-means*) is another version of our simple method in which our semi-supervised clustering is substituted by k-means clustering. *JEaC(init D)* and *JEaC(init R)* correspond to optimizing Eq. (4) with respectively initializing D using Eq. (2) and initializing R by IEaC. For our methods, average and standard deviation of different runs are reported. As it can be seen, the initialization done by IEaC has critical effect on the performance. This can be justified by noting the information from structure of unlabeled data is leveraged when initializing R while such information is absent in initializing D .

For other methods, we use the results reported in their original publication. Note that some experimental settings

¹However we found out through experiment that increasing the number of clusters improves the accuracy.

Table 1: Accuracy score (%) of cluster assignments converted to labels using majority voting on ground truth labels on four zero-shot recognition benchmarks. Results are our method are average \pm std of three runs.

Clustering Method	Animals with Attributes	CUB-2011	aPascal-aYahoo	SUN Attribute
k-means	65.80	35.61	65.37	17.49
Proposed Clustering	70.74\pm0.32	42.63\pm0.07	69.93\pm3.4	45.50\pm1.32

Table 2: Classification accuracy in % on four public datasets: Animals with Attributes, CUB-2011, aPascal-aYahoo and SUN in form of average \pm std.

Feature	Method	Animals with Attributes	CUB-2011	aPascal-aYahoo	SUN
Shallow	(Li and Guo 2015)	38.2 \pm 2.3			18.9 \pm 2.5
	(Schuurmans and Tg 2015)	40.05 \pm 2.25		24.71 \pm 3.19	
	(Jayaraman and Grauman 2014)	43.01 \pm 0.07		26.02 \pm 0.05	56.18 \pm 0.27
GoogleNet	(Akata et al. 2015b)	66.7	50.1		
	(Changpinyo et al. 2016)	72.9	54.5		
	(Xian et al. 2016)	71.9	45.5		62.7
VGG-19	(Kodirov et al. 2015)	73.2	39.5	26.5	
	(Akata et al. 2015b)	61.9	50.1		
	(Zhang and Saligrama 2015b)	76.33 \pm 0.53	30.41 \pm 0.20	46.23 \pm 0.53	82.50 \pm 1.32
	(Zhang and Saligrama 2015a)	80.46 \pm 0.53	42.11 \pm 0.55	50.35 \pm 2.97	83.83 \pm 0.29
	IEaC (k-means)	86.34 \pm 0.13	52.48 \pm 0.60	48.03 \pm 1.56	75.75 \pm 1.06
	IEaC (Proposed Clustering)	86.38 \pm 0.56	53.10 \pm 0.43	48.00 \pm 0.69	80.66pm0.76
	JEaC (init D)	83.03	57.55	42.62	72.50
	JEaC (init R)	88.64\pm0.04	58.80\pm0.64	49.77 \pm 2.02	86.16\pm0.57

of these works may differ from those of ours. We did not reimplement any of the other methods and if the original paper does not report results on a data set we leave the corresponding cell as blank. Our method performs the best on three out of the four datasets (outperforms the others on all except to the aPascal-aYahoo dataset). This can be explained by the nature of the dataset in which class signatures obtained by averaging instance attributes are very similar. We suppose trying to learn more discriminative signatures from data can potentially improve the result. We investigate this in our future work.

The effectiveness of different components of our methods in further illustrated in Figure 2. As it can be seen in Figure 2c merely using mapping from Eq. (2) results in poor signature embeddings where domain shift problem is visible. However using the compatibility function based on cluster assignments, although the there is no change in mappings, label assignments are improved, our clustering (Figure 2e) performing better than k-means (Figure 2d). Finally using mapping from JEaC, the domain shift problem is substantially alleviated and far better signature embeddings are achieved.

Conclusion

In this paper, we proposed semi-supervised methods for zero-shot object recognition. We used the space of deep visual features as a semantic visual space and learned a linear transformation to map class signatures to this space such that the mapped signatures provide good representative of the corresponding instances. We utilized this property that the rich deep visual features provide a representation space in which samples of each class are usually condensed in a clus-

ter. In the proposed method that jointly learns the mapping of class signatures and the class assignments of unlabeled data, we used also unlabeled instances of unseen classes when learning the mapping to alleviate the domain shift problem. Experimental results showed that the proposed method generally outperformed the other recent methods.

References

- Akata, Z.; Perronnin, F.; Harchaoui, Z.; and Schmid, C. 2013. Label-embedding for attribute-based classification. In *Computer Vision and Pattern Recognition (CVPR), IEEE Conference on*, 819–826.
- Akata, Z.; Perronnin, F.; Harchaoui, Z.; and Schmid, C. 2015a. Label-embedding for image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* PP(99).
- Akata, Z.; Reed, S.; Walter, D.; Lee, H.; and Schiele, B. 2015b. Evaluation of Output Embeddings for Fine-Grained Image Classification. In *Computer Vision and Pattern Recognition (CVPR), IEEE Conference on*.
- Akata, Z.; Malinowski, M.; Fritz, M.; and Schiele, B. 2016. Multi-Cue Zero-Shot Learning with Strong Supervision. *arXiv preprint arXiv:1603.08754*.
- Arthur, D., and Vassilvitskii, S. 2007. k-means++: the advantages of careful seeding. In *In Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 1027–1035.
- Ba, J.; Swersky, K.; Fidler, S.; and Salakhutdinov, R. 2015. Predicting Deep Zero-Shot Convolutional Neural Networks using Textual Descriptions. *arXiv preprint arXiv:1506.00511*.
- Changpinyo, S.; Chao, W.; Gong, B.; and Sha, F. 2016. Synthesized classifiers for zero-shot learning. *CoRR* abs/1603.00550.
- Chapelle, O.; Schölkopf, B.; and Zien, A. 2006. *Semi-Supervised Learning*. Cambridge, MA: MIT Press.

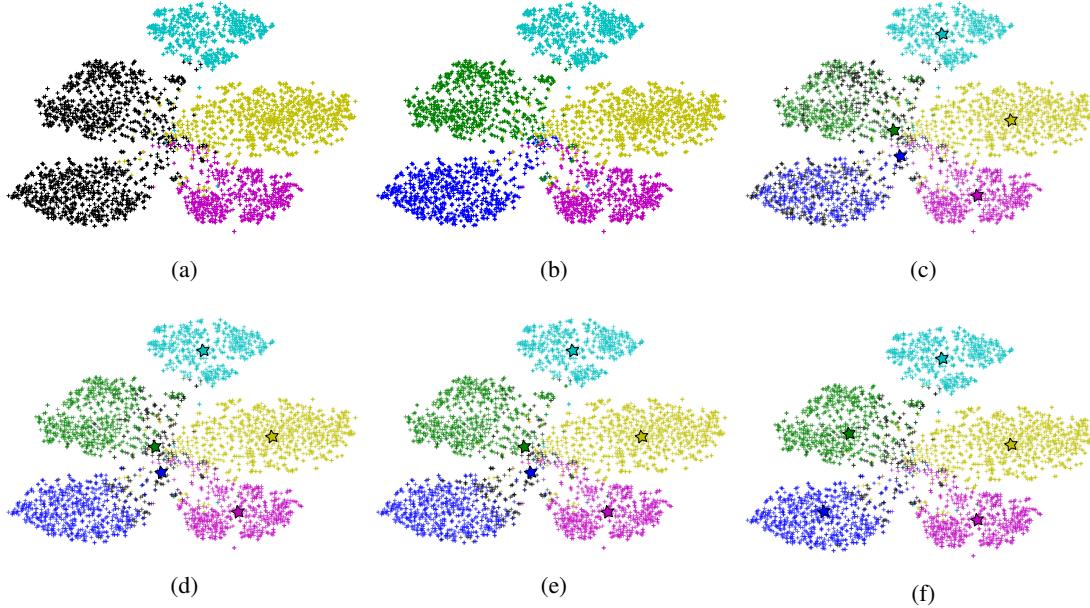


Figure 2: t-SNE embedding of five classes from AwA, three seen: antelope (magenta), grizzly bear (yellow), killer whale (cyan) and two unseen: chimpanzee (blue), giant panda (green). Images shown by plus signs and embedding of class signatures in images space by stars. in figures b-f black points denote assignment to a class other five classes shown here. **b)** Points colored according to their ground truth labels **c)** Signatures mapped to image spacing using Eq. (2). Then classification done using nearest neighbor **d)** Classification done by our compatibility function on cluster assignments from k-means **e)** Classification by our compatibility function using our supervised clustering **f)** Class signatures mapping and cluster assignment by JEaC

- Elhoseiny, M.; Saleh, B.; and Elgammal, A. 2013. Write a classifier: Zero-shot learning using purely textual descriptions. In *Computer Vision (ICCV), IEEE Conference on*, 2584–2591.
- Farhadi, A.; Endres, I.; Hoiem, D.; and Forsyth, D. 2009. Describing Objects by Their Attributes. In *Computer Vision and Pattern Recognition (CVPR), IEEE Conference on*, 1778–1785.
- Frome, A.; Corrado, G. S.; Shlens, J.; Bengio, S.; Dean, J.; Ranzato, M.; and Mikolov, T. 2013. DeViSE: A Deep Visual-Semantic Embedding Model. In *Advances in Neural Information Processing Systems (NIPS) 26*, 2121–2129.
- Fu, Y., and Sigal, L. 2016. Semi-supervised Vocabulary-informed Learning. *arXiv preprint arXiv:1604.07093*.
- Fu, Y.; Hospedales, T. M.; Xiang, T.; Fu, Z.; and Gong, S. 2014. Transductive multi-view embedding for zero-shot recognition and annotation. In *Computer Vision (ECCV), European Conference on*, volume 6315.
- Hoiem, D.; Divvala, S. K.; and Hays, J. H. 2008. Pascal voc 2008 challenge.
- Jayaraman, D., and Grauman, K. 2014. Zero-shot recognition with unreliable attributes. In *Advances in Neural Information Processing Systems (NIPS) 27*. 3464–3472.
- Kodirov, E.; Xiang, T.; Fu, Z.; and Gong, S. 2015. Unsupervised Domain Adaptation for Zero-Shot Learning. In *Computer Vision (ICCV), IEEE Conference on*, 2927–2936.
- Lampert, C.; Nickisch, H.; and Harmeling, S. 2009. Learning to detect unseen object classes by between-class attribute transfer. In *Computer Vision and Pattern Recognition (CVPR), IEEE Conference on*, 951–958.
- Larochelle, H.; Erhan, D.; and Bengio, Y. 2008. Zero-data learn-
ing of new tasks. In *National Conference on Artificial Intelligence (AAAI)*, 646–651.
- Li, X., and Guo, Y. 2015. Max-margin zero-shot learning for multi-class classification. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, 626–634.
- Mahajan, D.; Sellamanickam, S.; and Nair, V. 2011. A joint learning framework for attribute models and object descriptions. In *Computer Vision (ICCV), IEEE International Conference on*, 1227–1234.
- Norouzi, M.; Mikolov, T.; Bengio, S.; Singer, Y.; Shlens, J.; Frome, A.; Corrado, G.; and Dean, J. 2014. Zero-shot learning by convex combination of semantic embeddings. In *International Conference on Learning Representations*.
- Palatucci, M.; Hinton, G.; Pomerleau, D.; and Mitchell, T. M. 2009. Zero-shot learning with semantic output codes. In *Advances in Neural Information Processing Systems (NIPS) 22*. 1410–1418.
- Patterson, G.; Xu, C.; Su, H.; and Hays, J. 2014. The sun attribute database: Beyond categories for deeper scene understanding. *International Journal of Computer Vision* 108(1-2):59–81.
- Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.; Cournapeau, D.; Brucher, M.; Perrot, M.; and Duchesnay, E. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12:2825–2830.
- Romera-Paredes, B., and Torr, P. H. S. 2015. An Embarrassingly Simple Approach to Zero-shot Learning. *Journal of Machine Learning Research* 37.
- Schuurmans, D., and Tg, A. B. 2015. Semi-Supervised Zero-Shot

Classification with Label Representation Learning. In *Computer Vision (ICCV), IEEE Conference on*.

Scott Reed, Zeynep Akata, Honglak Lee, B. S. 2016. Learning Deep Representations of Fine-Grained Visual Descriptions. *CVPR*.

Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *CoRR*.

Suzuki, M.; Sato, H.; Oyama, S.; and Kurihara, M. 2014. Transfer learning based on the observation probability of each attribute. In *Systems, Man and Cybernetics (SMC), IEEE International Conference on*, 3627–3631.

Wah, C.; Branson, S.; Welinder, P.; Perona, P.; and Belongie, S. 2011. The Caltech-UCSD Birds-200-2011 Dataset. Technical report.

Wang, X., and Ji, Q. 2013. A unified probabilistic approach modeling relationships between attributes and objects. In *Computer Vision (ICCV), IEEE International Conference on*, 2120–2127.

Xian, Y.; Akata, Z.; Sharma, G.; Nguyen, Q.; Hein, M.; and Schiele, B. 2016. Latent Embeddings for Zero-shot Classification. *arXiv preprint arXiv:1603.08895*.

Yu, X., and Aloimonos, Y. 2010. Attribute-based transfer learning for object categorization with zero/one training example. In *Computer Vision (ECCV), European Conference on*, volume 6315. 127–140.

Yu, F. X.; Cao, L.; Feris, R. S.; Smith, J. R.; and Chang, S.-F. 2013. Designing Category-Level Attributes for Discriminative Visual Recognition. In *Computer Vision and Pattern Recognition (CVPR), IEEE Conference on*, 771–778.

Zhang, Z., and Saligrama, V. 2015a. Zero-shot learning via joint latent similarity embedding. *arXiv preprint arXiv:1511.04512*.

Zhang, Z., and Saligrama, V. 2015b. Zero-Shot Learning via Semantic Similarity Embedding. In *Computer Vision (ICCV), IEEE Conference on*.