

2016

# Learning joint feature adaptation for zero-shot recognition

---

Ziming Zhang, Venkatesh Saligrama. 2016. "Learning Joint Feature Adaptation for Zero-Shot Recognition." arXiv preprint arXiv:1611.07593,

<https://hdl.handle.net/2144/29425>

*"Downloaded from OpenBU. Boston University's institutional repository."*

# Learning Joint Feature Adaptation for Zero-Shot Recognition

Ziming Zhang

Mitsubishi Electric Research Laboratories  
201 Broadway, Cambridge, MA 02139-1955  
zzhang@merl.com

Venkatesh Saligrama

ECE, Boston University  
8 Saint Mary's Street, Boston, MA 02215  
srv@bu.edu

## Abstract

Zero-shot recognition (ZSR) aims to recognize target-domain data instances of unseen classes based on the models learned from associated pairs of seen-class source and target domain data. One of the key challenges in ZSR is the relative scarcity of source-domain features (e.g. one feature vector per class), which do not fully account for wide variability in target-domain instances.

In this paper we propose a novel framework of learning data-dependent feature transforms for scoring similarity between an arbitrary pair of source and target data instances to account for the wide variability in target domain. Our proposed approach is based on optimizing over a parameterized family of local feature displacements that maximize the source-target adaptive similarity functions. Accordingly we propose formulating zero-shot learning (ZSL) using latent structural SVMs to learn our similarity functions from training data. As demonstration we design a specific algorithm under the proposed framework involving bilinear similarity functions and regularized least squares as penalties for feature displacement. We test our approach on several benchmark datasets for ZSR and show significant improvement over the state-of-the-art. For instance, on aP&Y dataset we can achieve 80.89% in terms of recognition accuracy, outperforming the state-of-the-art by 11.15%.

## 1. Introduction

While there has been significant progress on supervised large-scale classification in recent years [43], the lack of sufficient annotated training data uniformly across all classes [8, 4] has been a bottleneck in achieving acceptable performance. At a basic level, in these cases we encounter situations where we may have sufficient annotated training data for some of the classes and little or even no annotated data to train supervised classifiers for the other interesting classes. In this context, a fundamental question that arises is as to how to leverage training data for observed classes

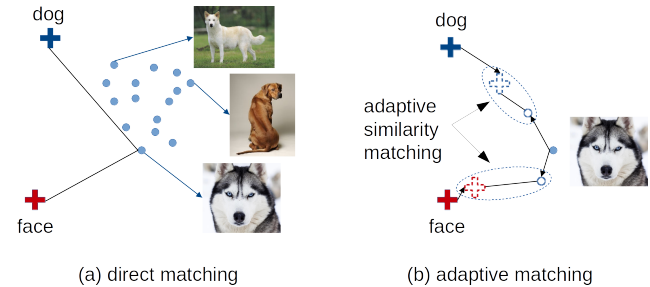


Figure 1. Illustration of our intuition behind learning joint data-dependent feature adaptation for ZSR. Here + and o denote source and target domain data embeddings, respectively, colors denote different classes, filled/empty shapes denote original/adapted feature embeddings. Using the direct matching in (a) the image will be mis-classified as “face” (based on distance measure), while using the adaptive matching in (b) it will be classified correctly as “dog”.

for recognition of rare or unobserved classes.

One possible scenario is when we have data from a different domain that can be collected easily, as is assumed to be the case in zero-shot recognition (ZSR). In ZSR we are given *source* and *target* domains as training data belonging to a sub-collection of classes, forming *seen* or *observed* classes. No training data is available for other *unseen* classes. The class information in source domain is described in a variety of ways such as attribute vectors [14, 28, 33, 36, 41], language words/phrases [6, 16, 45], or even learned classifiers [53]. Target domain is described by a joint distribution of data (e.g. images or videos) and labels [28, 50].

The challenge in ZSR lies in learning models based on seen-class data that can generalize to unseen classes. In this context, many perspectives for zero-shot learning (ZSL) have been proposed, including unseen classifier prediction from source domain data based on learning attribute classifiers [27], learning similarity functions between source and target domains to score similarity for unseen classes [1], and manifold embedding methods based on identifying inter-class relationships in source and target data that can be aligned during test time [56, 57, 11].

Nevertheless, the challenge posed by the relative sparseness of source domain descriptions in recognition has not

been fully considered. In particular, the target-domain data exhibits significant intra-class variation (*e.g.* appearance and poses). On the other hand the source domain information is relatively sparse and typically amounting to a single attribute vector. This is generally insufficient to account for all of the intra-class variation.

Fig. 1 illustrates this point. In the joint embedded feature space, the target-domain data distributions of certain classes (*e.g.* “dog” class in the figure) are relatively flat and consist of data instances with large variation. This issue leads us to view the presented source domain vector as a “mean-value” over all candidate (or alternative) source vectors. During test time for a given target instance we optimize the matching over all possible source and target candidates in a neighborhood of the presented source and target instances. During training we propose learning a data-dependent feature transform chosen from a parameterized family of displacement functions that maximizes the similarity between an arbitrary source and target instances. We learn our similarity functions from training data using *latent structural SVMs*. As demonstration we design a specific algorithm under the proposed framework involving bilinear similarity functions and regularized least squares as penalties for displacement.

To illustrate how this would work, consider again Fig. 1. In test time our proposed approach manifests as new features (*i.e.* empty “+” in the figure) that adapt to the potential contents in target data instance. This leads to significantly richer representations than the provided source-domain vectors. Our proposed approach also induces displacements in the target domain data instances. These displaced features (*i.e.* empty “o” in the figure) in turn adapt to source domain features. This process is akin to de-noising of presented target-domain data/features. As we see, by using the new features, the dog face image is correctly classified based on similarity measure between the new data-dependent adapted features, as illustrated in Fig. 1(b).

**Contributions:** In this paper we introduce a novel *adaptive similarity function* for comparing an arbitrary pair of source and target domain data instances. This function, in test time and adaptively in a data-dependent way, determines the similarity between presented source and target instances.

We propose considering optimizing over a parameterized bilinear family of functions for our cross-domain similarity measure. Alternating optimization is utilized to efficiently estimate (globally) best adapted features within a constrained family of displacements. In this context we show that the compatibility function defined in [1, 2] is indeed a special case of our similarity function.

To learn the parameters in our adaptive similarity function, we further propose formulating the ZSL problem using *latent structural SVMs*. The latent part comes from the adapted (latent) features, which are considered as the latent variables in the formulation. The structural part arises from

the structures of label embeddings as did in [1, 2].

We test our approach on four benchmark image datasets for ZSL/ZSR, namely, aP&Y, AwA, CUB and SUN-attribute. Under both standard and transductive settings, our approach outperforms the state-of-the-art significantly.

## 1.1. Related Work

In general ZSL/ZSR approaches can be divided into two categories: standard setting and transductive setting. Recently zero-shot approaches have been successfully applied to several visual tasks such as event detection [50, 10, 13], action recognition [20], and image tagging [55]. Below we primarily describe learning approaches in this context.

**Standard Setting:** In test time, the source-domain descriptors for unseen classes are all given at once. Our task is to sequentially recognize target-domain instances as they are revealed *one at a time*.

In this context, several works in the literature are based on training attribute classifiers which directly map target-domain data into source-domain attribute space [35, 28, 31, 49, 53, 54, 32, 22, 42, 3]. The resulting attribute classifiers do not fully account for data noise in source (*e.g.* ambiguity or mislabeling in attributes) and target (*e.g.* large variation because of the changes of appearance, poses, *etc.*) domains.

Linear and nonlinear embedding approaches [1, 2, 16, 34, 45, 29, 30, 39, 5, 25, 56, 57, 11, 51, 9, 48] have attracted attention recently. The basic idea of these methods is to embed the source and target domain features into a Kronecker product embedding space. For instance, Akata *et al.* [1, 2] proposed label embedding to map class labels into a high dimensional vector space (*e.g.* source-domain attribute space), and measure cross-domain similarities using a bilinear function whose parameters are learned using structured SVMs. Zhang and Saligrama [57] proposed a joint learning framework to learn the latent embeddings for both domains and utilized them for similarity measure. Changpinyo *et al.* [11] proposed a learning method to generate synthesized classifiers for unseen classes. Bucher *et al.* [9] proposed a metric learning based formulation to improve semantic embedding consistency, achieving the best performance on the four benchmark datasets under the standard setting in the current literature, to our best knowledge. The underlying assumption behind such approaches is that there exist (hidden) corresponding matches between source-domain feature vectors and target-domain data distributions, *e.g.* one-to-one match [1, 2, 57] or one-to-many match [51]. In this context there are other related proposed methods such as semantic transfer propagation [40], random forest based approaches [23], semantic manifold distance [19] approaches, and similarity calibration method [12]. Nevertheless, the issue of source-domain sparsity and the resulting imbalance with target-domain data is not fully accounted for in these methods.

Our proposed method explicitly focuses on handling the scarcity issue of source-domain data by learning data-dependent latent features. This in turn accounts for the large data variation in target domain implicitly so that the cross-domain matches can be improved.

**Transductive Setting:** Recently researchers have begun to incorporate test-time unseen-class data in target domain into ZSL/ZSR as unlabeled data analogous to the transductive setting. This has led to approaches that attempt to account for domain shift [25, 17, 18, 21, 58]. In this setting, during test time, we are given a list of all unlabelled target instances in addition to unseen-class source-domain descriptions. Potentially these methods can be used in conjunction with any similarity learning procedure trained on seen-class data, as demonstrated in [58], to score similarity between unseen classes and target domain data instances.

While much of the focus of this paper is on the standard setting, in our experimental section we also test our learning algorithm in the transductive mode to benchmark our performance in the transductive setting.

## 2. Our Approach

### 2.1. ZSL/ZSR Problem Setup

In the training stage, we are given a set of observed classes  $\mathcal{L}_o$ . For source domain, attribute vectors (or label embeddings) in the form of  $\{\psi(y)\}$ ,  $\forall y \in \mathcal{L}_o$ , are provided. Typically there exists only one vector per class. Corresponding target domain data instances  $x \in \mathcal{X}$  and feature embeddings  $\phi(x)$  associated with the observed source labels are also provided for training. We aggregate training data as  $\mathbb{O} = \{(\phi(x_i), \psi(y_i), y_i), \forall i \in \mathbb{T}\}$ , where  $i$  denotes the  $i$ -th training data instance in target domain.

Our goal is to learn a prediction model, by leveraging observed training data,  $\mathbb{O}$ , such that it generalizes well to unobserved data instances and classes during test time.

In the testing stage, a set of source vectors corresponding to unobserved classes  $\mathcal{L}_u$  are revealed. For a given unobserved data instance  $\bar{x}$  from target domain, the task is to identify the source vector among those unobserved classes that corresponds to  $\bar{x}$ . Abstractly, our decision rule is based on maximizing a posterior probability (MAP) conditioned on all the available data:

$$\bar{y}^* = \arg \max_{\bar{y} \in \mathcal{L}_u} \mathbb{P}_{\psi, \phi}(\bar{y} | \bar{x}; \mathbb{O}), \quad (1)$$

where  $\mathbb{P}_{\psi, \phi}(\cdot)$  denotes the posterior probability tuned to the embedding functions  $\psi, \phi$ . In what follows we drop the parameter dependence on  $\psi, \phi$  for notational simplicity, since we assume that these embedding functions are provided a priori. The posterior probability is unknown and must be learned from training data. We describe our proposed approach in the following section.

## 2.2. General Learning Framework

### 2.2.1 Parameterized Family of Posterior Distributions

We face two fundamental challenges in ZSR.

First, target instances and labels for unobserved classes are not known during training. Therefore, proposed methods must base its recognition on scoring the similarity between an arbitrary source descriptor and a target instance.

Second, source vectors in ZSR are sparse and typically we only observe a single source vector per class. On the other hand there is significant variability in the target domain. Consequently, the source vectors serve only as “average” attribute descriptors across the target domain instances. The source descriptor that best matches a target instance is a vector that is typically close to but not necessarily equal to the given source domain vector. We propose optimizing over all such vectors in both learning and test time to determine the optimal matching source descriptors.

In this context we propose a family of posterior distributions: To account for relative sparseness of source domain descriptors and large variability of target domain instances we introduce new data-dependent feature vectors  $\mathbf{z}_s, \mathbf{z}_t$  corresponding to source and target domains, respectively. To ensure that these feature vectors are “typically” close to the given source and target data pair we introduce a displacement penalty term  $d_\omega(x, y, \mathbf{z}_s, \mathbf{z}_t)$  parametrized by  $\omega$ . To score similarity between source and target domain data we propose a scoring function,  $s_{\mathbf{W}}(\mathbf{z}_s, \mathbf{z}_t)$ , that scores similarity between the new data-dependent feature vectors parametrized by a matrix  $\mathbf{W}$ . This leads to the following posterior probability:

$$\mathbb{P}(y, \mathbf{z}_s, \mathbf{z}_t | x; \mathbf{W}, \omega) \propto \exp(s_{\mathbf{W}}(\mathbf{z}_s, \mathbf{z}_t) - d_\omega(x, y, \mathbf{z}_s, \mathbf{z}_t)). \quad (2)$$

In order to compute the posterior  $\mathbb{P}(y|x)$  we can marginalize  $\mathbb{P}(y, \mathbf{z}_s, \mathbf{z}_t | x; \mathbf{W}, \omega)$  over variables  $\mathbf{z}_s \in \mathcal{Z}_s, \mathbf{z}_t \in \mathcal{Z}_t$ , where  $\mathcal{Z}_s, \mathcal{Z}_t$  denote their corresponding feasible domains (e.g. simplex). However, in general this calculation will be very difficult given arbitrary parameter spaces, and typically Bayesian parametrization is often involved (e.g. [38]) to simplify the calculation. Alternatively the posterior can be upper-bounded by the maximum value over the variables, as did in [57], which can be very computationally efficient and demonstrated with good performance for ZSR as well.

Therefore, here we adopt the strategy in [57] and take the maximum for posterior approximation purpose. This leads naturally to our *adaptive similarity function* as below for scoring each target data instance with a class label:

$$\begin{aligned} f(x, y; \mathbf{W}, \omega) &\triangleq \max_{\mathbf{z}_s \in \mathcal{Z}_s, \mathbf{z}_t \in \mathcal{Z}_t} \mathbb{P}(y, \mathbf{z}_s, \mathbf{z}_t | x; \mathbf{W}, \omega) \\ &= \max_{\mathbf{z}_s \in \mathcal{Z}_s, \mathbf{z}_t \in \mathcal{Z}_t} \left\{ s_{\mathbf{W}}(\mathbf{z}_s, \mathbf{z}_t) - d_\omega(x, y, \mathbf{z}_s, \mathbf{z}_t) \right\}. \quad (3) \end{aligned}$$

Intuitively our similarity function allows the features to move from their original locations in the feature space (*i.e.* adaptation) to achieve a higher similarity score within a neighborhood (feature displacements incur penalties). Our function in Eq. 3 thus attempts to achieve a balance between these two objectives. In fact similar strategy has been widely used in deformable part models (DPM) [15], where 2D locations for parts are considered as adapted features.

### 2.2.2 Learning with Latent Structural SVMs

The parameters  $\mathbf{z}_s, \mathbf{z}_t$  in Eq. 3 play the role of latent variables for given values of  $\mathbf{W}, \omega$ . Consequently, we can pose the problem as a latent structural SVM problem by viewing the label variable  $y$  as taking values from a structured output space:

$$\begin{aligned} \min_{\mathbf{W} \in \mathcal{W}, \omega \in \Omega, \xi} \mathcal{R}_1(\mathbf{W}) + \mathcal{R}_2(\omega) + \sum_i \xi_i \quad (4) \\ \text{s.t. } f(x_i, y_i; \mathbf{W}, \omega) - f(x_i, y; \mathbf{W}, \omega) \geq \Delta(y_i, y) - \xi_i, \\ \forall i, \xi_i \geq 0, x_i \in \mathcal{X}, y_i, y \in \mathcal{L}_o, \end{aligned}$$

where  $\mathcal{R}_1, \mathcal{R}_2$  denote two regularization functions (*e.g.*  $\ell_2$ -norm regularizers) for parameters  $\mathbf{W}, \omega$ , respectively,  $\Delta$  denotes a penalty term measuring the difference between the ground-truth label  $y_i$  and an arbitrary label  $y$ ,  $\mathcal{W}, \Omega$  denote the feasible domains for  $\mathbf{W}, \omega$ , respectively, and  $\xi_i, \forall i$  is a slack variable. The cutting-plane algorithm [52] can be used for general training purpose.

In test time, we replace the probability term in Eq. 1 with our adaptive similarity function in Eq. 3 to rewrite the decision rule for ZSR as follows:

$$\bar{y}^* = \arg \max_{\bar{y} \in \mathcal{L}_u} f(\bar{x}, \bar{y}; \mathbf{W}, \omega), \quad \forall \bar{x}. \quad (5)$$

### 2.3. Bilinear Adaptive Similarity Functions

For the purpose of demonstration we describe one instance of an adaptive similarity function that can be utilized in our general learning framework.

Specifically we design the similarity term  $s_{\mathbf{W}}(\mathbf{z}_s, \mathbf{z}_t)$  in Eq. 3 as a bilinear function. These type of functions have been widely used in recent ZSL literature, *e.g.* [1, 2, 57], and has been shown to achieve state-of-the-art performance. For the penalty term, we simply adopt the regularized least square loss for the displacement. Putting these together, we propose the following adaptive similarity function:

$$\begin{aligned} f(x, y; \mathbf{W}, \omega) = \max_{\mathbf{z}_s \in \mathcal{Z}_s, \mathbf{z}_t \in \mathcal{Z}_t} \left\{ \mathbf{z}_t^T \mathbf{W} \mathbf{z}_s - \frac{\omega_1}{2} \|\mathbf{z}_t - \phi(x)\|_2^2 \right. \\ \left. - \frac{\omega_2}{2} \|\mathbf{z}_s - \psi(y)\|_2^2 - \frac{\omega_3}{2} \|\mathbf{z}_t\|_2^2 - \frac{\omega_4}{2} \|\mathbf{z}_s\|_2^2 \right\}, \quad (6) \end{aligned}$$

where  $\mathbf{W} \in \mathbb{R}^{d_t \times d_s}$  is a weighting matrix between  $\mathbf{z}_t \in \mathbb{R}^{d_t}$  and  $\mathbf{z}_s \in \mathbb{R}^{d_s}$ ,  $\omega = [\omega_1; \omega_2; \omega_3; \omega_4]$  is a 4D vector

controlling the trade-off between similarity and penalty, and  $\|\cdot\|_2$  denotes the  $\ell_2$  norm of a vector. In general we can utilize *alternating optimization* (AO) to solve Eq. 6 as follows:

$$\mathbf{z}_t = \arg \min_{\mathbf{z} \in \mathcal{Z}_t} \left\{ \omega_{13} \left\| \mathbf{z} - \left( \frac{\omega_1 \phi(x)}{\omega_{13}} + \frac{\mathbf{W} \mathbf{z}_s}{\omega_{13}} \right) \right\|_2^2 \right\}, \quad (7)$$

$$\mathbf{z}_s = \arg \min_{\mathbf{z} \in \mathcal{Z}_s} \left\{ \omega_{24} \left\| \mathbf{z} - \left( \frac{\omega_2 \psi(y)}{\omega_{24}} + \frac{\mathbf{z}_t^T \mathbf{W}}{\omega_{24}} \right) \right\|_2^2 \right\}, \quad (8)$$

where  $\omega_{13} = \omega_1 + \omega_3$  and  $\omega_{24} = \omega_2 + \omega_4$ .

Ideally, we would like to have a decision rule that during test time using Eq. 6 converges to a (unique) global solution<sup>1</sup> for an arbitrary pair of source and target data instances. This is because we can then be certain that the similarity scores are unique and reliable. Therefore, below we provide some general and useful properties about of the similarity function in Eq. 6.

**Property 1** (Global Optimality). *Let us define a new matrix*

$$\mathbf{H} = \begin{bmatrix} \omega_{13} \mathbf{I}_{d_t \times d_t} & -\mathbf{W} \\ -\mathbf{W}^T & \omega_{24} \mathbf{I}_{d_s \times d_s} \end{bmatrix}, \quad (9)$$

where  $\mathbf{I}_{d_t \times d_t}$  and  $\mathbf{I}_{d_s \times d_s}$  denote two identity matrices with sizes of  $d_t \times d_t$  and  $d_s \times d_s$  entries, respectively. Then if  $\mathbf{H}$  is positive definite (PD) and  $\mathcal{Z}_s, \mathcal{Z}_t$  are nonempty closed convex sets, there exists a unique global solution for Eq. 6.

*Proof.* Eq. 6 can be rewritten with  $\mathbf{H}$  in Eq. 9 as follows:

$$\begin{aligned} f(x, y; \mathbf{W}, \omega) \\ = \min_{\mathbf{z}} \left\{ \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z} - \mathbf{z}^T g(x, y; \omega) + h(x, y, \omega) \right\}, \quad (10) \end{aligned}$$

where  $\mathbf{z} = \begin{bmatrix} \mathbf{z}_t \\ \mathbf{z}_s \end{bmatrix}$ ,  $g(x, y; \omega) = \begin{bmatrix} \omega_1 \phi(x) \\ \omega_2 \psi(y) \end{bmatrix}$ , and  $h(x, y, \omega) = \frac{\omega_1}{2} \|\phi(x)\|_2^2 + \frac{\omega_2}{2} \|\psi(y)\|_2^2$ . Since  $\mathcal{Z}_s, \mathcal{Z}_t$  are nonempty closed convex sets, the feasible domain for  $\mathbf{z}$  is nonempty closed convex as well. Based on [7], we can prove this property.  $\square$

**Property 2** (Global Convergence of AO). *Under the conditions in Property 1, the alternating optimization in Eq. 7 and Eq. 8 can guarantee global convergence.*

*Proof.* Due to matrix  $\mathbf{H}$  being PD, we can have  $\omega_{13} > 0$  and  $\omega_{24} > 0$ . Further since  $\mathcal{Z}_s, \mathcal{Z}_t$  are nonempty closed convex sets, both Eq. 7 and Eq. 8 define convex optimization problems (see [7]), respectively. Now based on Property 1 we can prove this property.  $\square$

<sup>1</sup>In this paper we only focus on utilizing convex optimization to achieve such global solutions, although minimizing concave functions over nonempty closed convex sets may result in global solutions as well, but it is much harder to be solved and, more importantly, without any guarantee.

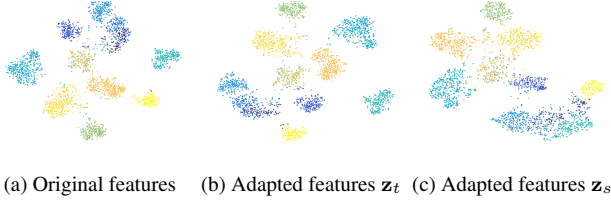


Figure 2. t-SNE visualization on aP&Y between (a) target-domain original features, (b) target-domain adapted features  $\mathbf{z}_t$  when matching with an unseen class, and (c) corresponding source-domain adapted features  $\mathbf{z}_s$ . Here each color represents a unique unseen class across the sub-figures.

**Property 3** (Local Convergence of AO). *If  $\omega_{13} > 0$ ,  $\omega_{24} > 0$ , and  $\mathcal{Z}_s, \mathcal{Z}_t$  are nonempty closed convex sets, then the alternating optimization in Eq. 7 and Eq. 8 can guarantee to converge to local optima.*

*Proof.* Please refer to the proof for Property 2.  $\square$

**Property 4** (Extreme Case). *Suppose that all the vectors and matrix in Eq. 6 are upper-bounded. Then we have*

$$\lim_{\omega_1, \omega_2 \rightarrow +\infty, \omega_3, \omega_4 \rightarrow 0} f(x, y; \mathbf{W}, \boldsymbol{\omega}) = \phi(x)^T \mathbf{W} \psi(y). \quad (11)$$

From Property 4 we can easily see that our adaptive similarity function in Eq. 6 can be taken as the generalization of the bilinear compatibility function defined in [1, 2], and so does our learning framework in Eq. 4 accordingly.

## 2.4. A Specific Learning Algorithm

With various feasible domains  $\mathcal{Z}_s, \mathcal{Z}_t$ , we can design different adaptive similarity functions accordingly. Particularly here we define

$$\mathcal{Z}_s = \{\mathbf{z}_s \mid \|\mathbf{z}_s\|_2^2 \leq \gamma_s, \gamma_s \geq 0, \forall \mathbf{z}_s\}, \quad (12)$$

$$\mathcal{Z}_t = \{\mathbf{z}_t \mid \|\mathbf{z}_t\|_2^2 \leq \gamma_t, \gamma_t \geq 0, \forall \mathbf{z}_t\}. \quad (13)$$

That is, we define  $\mathcal{Z}_s, \mathcal{Z}_t$  to be sufficiently large sets which contain *any possible* source or target adapted feature embedding, respectively<sup>2</sup>. Our reasoning for this choice is its simplicity and our need for high computational efficiency.

Then by setting the first derivative of  $f$  over  $\mathbf{z}$  to 0, *i.e.*  $\frac{\partial f}{\partial \mathbf{z}} = 0$ , we can easily get the close-form solution for  $\mathbf{z}$ , equivalently for  $\mathbf{z}_s, \mathbf{z}_t$ , as follows:

$$\mathbf{z} = \mathbf{H}^\dagger g(x, y; \boldsymbol{\omega}), \quad (14)$$

where  $\dagger$  denotes the pseudo-inverse operation.

**Discussion:** Eq. 14 suggests a linear transform function of combining source and target information to generate the adapted features. Since  $\mathbf{H}$  is PD (and thus so is  $\mathbf{H}^\dagger$ ), the target domain data structures are fully preserved in  $\mathbf{z}_t$  while

<sup>2</sup>Intuitively we can set  $\gamma_s, \gamma_t \rightarrow +\infty$ , *i.e.* very large real numbers.

matching with a single source domain vector. Correspondingly the target data structures will have a larger impact in generating  $\mathbf{z}_s$  as well. Fig. 2 illustrates the distributions of different features using the test data in aP&Y dataset, which conform with our analysis.

Next we substitute Eq. 6, 12 and 13 into Eq. 4 to learn the parameters in  $f$ . Note that in order to achieve global optimality in Property 1, the learned parameters must guarantee that matrix  $\mathbf{H}$  in Eq. 9 is PD. This leads us to the following learning problem:

$$\begin{aligned} \min_{\mathbf{W}, \boldsymbol{\omega}, \boldsymbol{\xi}} \quad & \frac{\lambda_1}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_2}{2} \|\boldsymbol{\omega}\|_2^2 + \sum_i \xi_i \quad (15) \\ \text{s.t.} \quad & f(x_i, y_i; \mathbf{W}, \boldsymbol{\omega}) - f(x_i, y; \mathbf{W}, \boldsymbol{\omega}) \geq \Delta(y_i, y) - \xi_i, \\ & \mathbf{H}(\mathbf{W}, \boldsymbol{\omega}) \succ \mathbf{0}, \\ & \forall i, \xi_i \geq 0, x_i \in \mathcal{X}, y_i, y \in \mathcal{L}_o, \end{aligned}$$

where  $\mathbf{H}(\mathbf{W}, \boldsymbol{\omega}) \equiv \mathbf{H}$  in Eq. 9 with  $\mathbf{W}, \boldsymbol{\omega}$  as parameters, “ $\succ \mathbf{0}$ ” denotes the PD constraint which makes it very difficult to solve the problem, and  $\lambda_1 \geq 0, \lambda_2 \geq 0$  are two predefined regularization parameters.

As a relaxation we tried to solve Eq. 15 without considering the PD constraint. However, we observed empirically that the learned parameters do not always satisfy the PD constraint using AO procedure. This leads to poor recognition performance. On the other hand, if we assume that the maximum  $\ell_1$  norm of the row vectors  $\mathbf{W}_{i,\cdot}, \forall i$  or column vectors  $\mathbf{W}_{\cdot,j}, \forall j$  in matrix  $\mathbf{W}$ , denoted by

$$\delta_W = \max \left\{ \max_i \|\mathbf{W}_{i,\cdot}\|_1, \max_j \|\mathbf{W}_{\cdot,j}\|_1 \right\}, \quad (16)$$

is non-zero and upper-bounded (which is always the case), we can obtain global optimality at least by manually setting parameter  $\boldsymbol{\omega}$  so that  $\omega_{13} \geq \delta_W$  and  $\omega_{24} \geq \delta_W$ . This creates a diagonally dominant matrix for  $\mathbf{H}$  and guarantees that PD is satisfied.

Based on this consideration, we chose not to learn parameter  $\boldsymbol{\omega}$  but instead set it manually to guarantee the PD constraint during training. We thus only learn parameter  $\mathbf{W}$ . Our learning formulation can now be rewritten as follows:

$$\begin{aligned} \min_{\mathbf{W}, \boldsymbol{\xi}} \quad & \frac{\lambda}{2} \|\mathbf{W}\|_F^2 + \sum_i \xi_i \quad (17) \\ \text{s.t.} \quad & f(x_i, y_i; \mathbf{W}, \boldsymbol{\omega}^*) - f(x_i, y; \mathbf{W}, \boldsymbol{\omega}^*) \geq \Delta(y_i, y) - \xi_i, \\ & \forall i, \xi_i \geq 0, x_i \in \mathcal{X}, y_i, y \in \mathcal{L}_o, \end{aligned}$$

where  $\boldsymbol{\omega}^*$  denotes the predefined parameter vector and  $\lambda \geq 0$  is a predefined constant. Since in our experiments the current ZSR problem is essentially equivalent to a multi-class prediction problem, we simply set  $\Delta(y_i, y) = 1$  if  $y_i \neq y$ , otherwise 0.

Table 1. Statistics of different benchmark image datasets.

Dataset	# instances	# attributes	# seen/unseen cls.
aP&Y	15,339	64 (continuous)	20 / 12
AwA	30,475	85 (continuous)	40 / 10
CUB-200-2011	11,788	312 (binary)	150 / 50
SUN Attribute	14,340	102 (binary)	707 / 10

In test time, by substituting Eq. 6, 9 and 14 into Eq. 5 we can rewrite the our decision function for ZSR as follows:

$$\begin{aligned} \bar{y}^* &= \arg \max_{\bar{y} \in \mathcal{L}_u} F(\bar{x}, \bar{y}; \mathbf{H}^\dagger, \omega^*) \\ &= \arg \max_{\bar{y} \in \mathcal{L}_u} \left\{ \frac{1}{2} g(\bar{x}, \bar{y}; \omega^*)^T \mathbf{H}^\dagger g(\bar{x}, \bar{y}; \omega^*) - h(\bar{x}, \bar{y}, \omega^*) \right\}. \end{aligned} \quad (18)$$

**Discussion:** Learning based on Eq. 17 has convergence issues due to the nature of latent structural SVMs. On the other hand, since our decisions are based on  $\mathbf{H}^\dagger$  explicitly as in Eq. 18, it might be possible to learn  $\mathbf{H}^\dagger$  approximately and efficiently by substituting similarity function  $F$  in Eq. 18 into structural SVMs [24]. It turns out that this learning strategy is equivalent to [1, 2] with source-domain feature augmentation, and thus leads to global convergence (under the multi-class prediction setting for ZSR). Empirically we tested this learning strategy and found marginal differences from [1, 2] in terms of recognition performance. Therefore we do not report these results in our experimental section.

### 3. Experiments

We follow the experimental settings in [57]. Specifically we test our method on four benchmark image datasets for zero-shot recognition, namely, aPascal & aYahoo (aP&Y) [14], Animals with Attributes (AwA) [26], Caltech-UCSD Birds-200-2011 (CUB-200-2011) [47], and SUN Attribute [37]. We summarize the statistics in each dataset and list them in Table 1.

For aP&Y, CUB-200-2011 and SUN Attribute datasets, we take the means of attribute vectors from the same classes to generate source domain data. For AwA dataset, we utilize the real-number attribute vectors since they are more discriminative. For all the datasets, we utilize MatConvNet [46] with the “imagenet-vgg-verydeep-19” pretrained model [44] to extract a 4096-dim CNN feature vector (*i.e.* the top layer hidden unit activations of the network) for each image (or bounding box). As suggested in [57] and with the same parameters we conduct dimension reduction for target-domain data and sparse coding for source-domain attribute vectors as well. All the predefined parameters in our method are tuned using cross-validation, similar to [56, 57]. We report our results averaged over 10 trials.

We utilize the same standard training/testing splits for zero-shot recognition on aP&Y and AwA as others, defined in these datasets. For CUB-200-2011, we follow [1] and

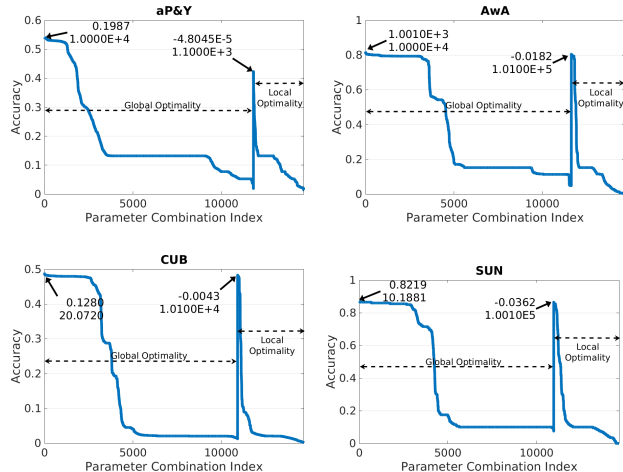


Figure 3. **Global optimality vs. local optimality:** Performance comparison over different parameter combinations of  $\omega$ . For both cases the reported results here have been sorted in a descending order. The numbers on the curves are the associated (top) smallest and (bottom) largest eigenvalues of matrix  $\mathbf{H}$  with the best results.

use the same 150 bird species as seen classes for training, and the other 50 species as unseen classes for testing. For SUN Attribute, we follow [23] and use the same 10 classes as unseen classes for testing (see their supplementary file), and the rest of them as seen classes for training.

On the four datasets, we evaluate the performance of the proposed approach in terms of recognition under the following two different settings:

**Standard vs. Transductive.** Our goal here is to benchmark the performance gains with different types of side-information in classification. Standard setting represents an extreme streaming situation. Transductive setting provides a less extreme scenario where during test-time we are given target instances all at once as in a batch-mode. The batch mode clearly provides information about the target data/feature distribution. The question arises as to how much we could benefit from this type of information. In this context we propose taking the similarities from our approach as inputs to the method proposed in [58] for recognition. Here we report our results averaged over 10 trials. In each trial we run our integrated approach for another 50 times and record the average as probabilities over unseen classes per target data. We predict class labels and report our performance based on this assignment probability matrix in each trial.

#### 3.1. Parameter Selection

In our method there are 5 predefined parameters, *i.e.*  $\lambda \geq 0$  and  $\omega = [\omega_1; \omega_2; \omega_3; \omega_4]$ . In this section, we will investigate the impact of  $\omega$  on recognition accuracy while fixing  $\lambda = 1$  without fine-tuning. Specifically we conduct a grid search over  $10^{-5:5}$  for each parameter in  $\omega$ , *i.e.* 11 choices per parameter and  $11^4 = 14,641$  param-

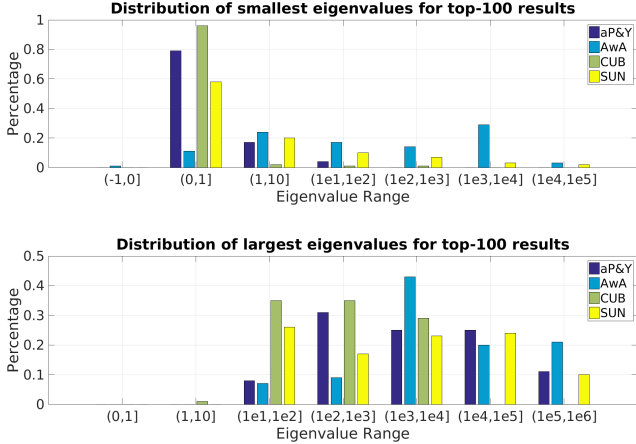


Figure 4. Illustration of distributions of (top) smallest and (bottom) largest eigenvalues for the corresponding top-100 results on each dataset.

ter combinations in total. In this way our adaptive similarity function in Eq. 6 always converges to a local maxima (*i.e.* local optimality) because of Property 3. Further when the corresponding matrix  $\mathbf{H}$  in Eq. 9 is PD, our function achieves global maxima (*i.e.* global optimality) because of Property 2. We utilize the smallest eigenvalue of  $\mathbf{H}$  as an indicator of being PD if it is positive, otherwise not.

First we investigate the effect of global/local optimality on test-time recognition accuracy, and illustrate the results in Fig. 3. Clearly we can observe that (1) *Quality*: The best global optimality results outperform their local optimality counterparts. (2) *Robustness*: The highest performing global optimality solution over different parameter combinations is more robust. This is indicated by the fact that for global optimality the curves for the top performance are going down very slowly, forming wide and relatively flat ranges, while for local optimality the curves are going down rapidly. The robustness here indeed suggests that with global optimality the parameter selection for  $\omega$  could have a sufficient number of choices, making it relatively easy. (3) *Eigenvalues*: In general the smallest eigenvalues tend to be close to 0, but the largest eigenvalues tend to be very large.

In order to accelerate the process of parameter selection, we provide some insights by looking at the distributions of both smallest and largest eigenvalues, as illustrated in Fig. 4. Generally speaking, smallest eigenvalues tend to concentrate on the range  $(0, 1]$ , while largest ones are uniformly distributed between  $10^1$  to  $10^6$  with slightly better focus on the range  $(10^3, 10^4]$ . This in turn suggests that for our method a good parameter combination for  $\omega$  may lead to a PD matrix  $\mathbf{H}$  (with high probability) whose smallest and largest eigenvalues lie in  $(0, 1]$  and  $(10^3, 10^4]$ , respectively. By quickly checking this condition, we can easily rule out most potential combinations. Further the big difference between the smallest and largest eigenvalues indicates a big difference between  $\omega_{13}$  and  $\omega_{24}$  as well. Moreover,

Table 2. Traditional zero-shot recognition accuracy comparison (%) with cited numbers in the form of mean $\pm$ standard deviation, grouped by the image feature types of (top) handcrafted features, (middle) other deep learning features, and (bottom) vgg-verydeep-19. Blanks indicate no reports for the datasets in the original papers.

Method	aP&Y	AwA	CUB	SUN
Farhadi <i>et al.</i> [14]	32.5			
Mahajan <i>et al.</i> [31]	37.93			
Wang and Ji [49]	45.05	42.78		
Rohrbach <i>et al.</i> [40]		42.7		
Yu <i>et al.</i> [53]		48.30		
Akata <i>et al.</i> [1]		43.5	18.0	
Mensink <i>et al.</i> [32]			14.4	
Lampert <i>et al.</i> [28]	19.1	40.5		52.50
J. and Grauman [23]	26.02 $\pm$ 0.05	43.01 $\pm$ 0.07		56.18 $\pm$ 0.27
R.-P. and Torr [42]	27.27 $\pm$ 1.62	49.30 $\pm$ 0.21		65.75 $\pm$ 0.51
Akata <i>et al.</i> [2]		66.7	<b>50.1</b>	
Qiao <i>et al.</i> [39]		66.46 $\pm$ 0.42	29.00 $\pm$ 0.28	
Changpinyo <i>et al.</i> [11]		72.9		
Xian <i>et al.</i> [51]		71.9	45.5	
Wang <i>et al.</i> [48]		75.99	33.48	
Lampert <i>et al.</i> [28]	38.16	57.23		72.00
R.-P. and Torr [42]	24.22 $\pm$ 2.89	75.32 $\pm$ 2.28		82.10 $\pm$ 0.32
SSE-INT [56]	44.15 $\pm$ 0.34	71.52 $\pm$ 0.79	30.19 $\pm$ 0.59	82.17 $\pm$ 0.76
SSE-ReLU [56]	46.23 $\pm$ 0.53	76.33 $\pm$ 0.83	30.41 $\pm$ 0.20	82.50 $\pm$ 1.32
SDL [57]	50.35 $\pm$ 2.97	79.12 $\pm$ 0.53	41.78 $\pm$ 0.52	83.83 $\pm$ 0.29
Bucher <i>et al.</i> [9]	<b>53.15<math>\pm</math>0.88</b>	77.32 $\pm$ 1.03	43.29 $\pm$ 0.38	<b>84.41<math>\pm</math>0.71</b>
<b>Ours: JFA</b>	52.04 $\pm$ 1.35	<b>81.03<math>\pm</math>0.88</b>	46.48 $\pm$ 1.67	84.10 $\pm$ 1.51

if  $\delta_W$  in Eq. 16 can be estimated, we may select parameter combinations more efficiently by constructing diagonal dominant matrices intentionally. This shows that the special structure of  $\mathbf{H}$  may provide important information to guide parameter selection.

### 3.2. Zero-Shot Recognition

In this section we compare our method, denoted by JFA, with other existing ZSL/ZSR approaches. This task is fundamentally about classification when a single target data instance is presented in test time.

#### 3.2.1 Standard Setting

Table 2 summarizes our comparison under the standard ZSR setting. Clearly ZSL methods leverage the advantages of deep learning features and achieve much better performance than those using handcrafted features. Different deep learning features can achieve comparable performance. Among all the competitors using vgg-verydeep-19 features, our method works the best, outperforming the state-of-the-art [9] by 1.34% on average. Notice that our experimental setting is exactly the same as [57], and in this case our method outperforms [57] significantly by 2.11%, but our standard deviations are slightly higher than [57]. The improvement comes from the nature of adaptive matching with better similarities, while the downside is mainly because our learning algorithm in Eq. 17 does not converge globally, leading to different local solutions even given the same training data.

To see this, let us take the CUB dataset for example. Initially CUB is created for fine-grained classification prob-

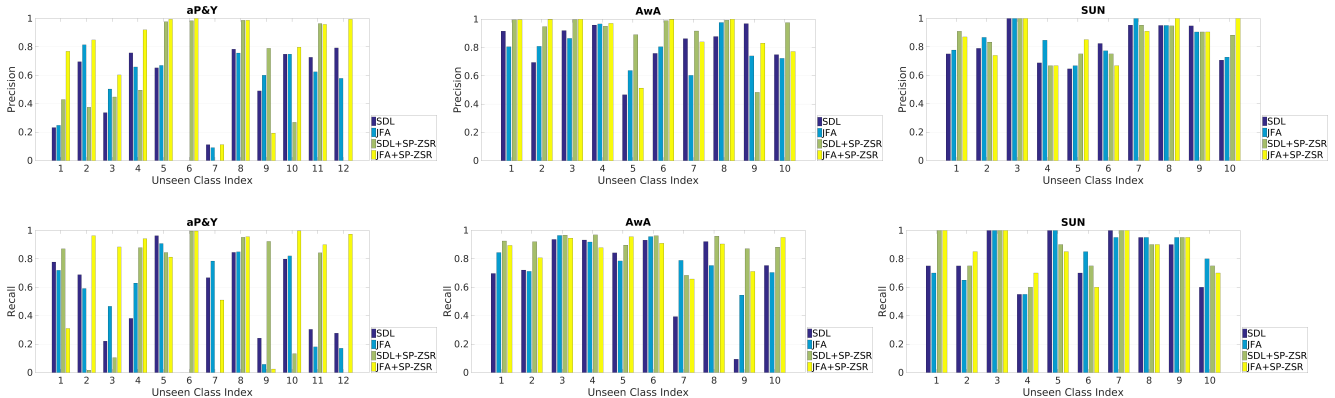


Figure 5. Performance comparison in terms of **(top)** precision and **(bottom)** recall among unseen classes indexed by the orders in the corresponding datasets. The numbers for the competitors are cited from [58].

Table 3. Transductive zero-shot recognition accuracy comparison (%) with cited numbers in the form of mean±standard deviation.

Method	aP&Y	AwA	CUB	SUN
Fu <i>et al.</i> [17]		47.1		
Fu <i>et al.</i> [18]		80.5	47.9	
Kodirov <i>et al.</i> [25]		75.6	40.2	
Guo <i>et al.</i> [21]	39.03	78.47		82.00
R.&T. [42]+SP-ZSR [58]	37.5	84.3		89.5
SDL [57]+SP-ZSR [58]	62.19±4.65	<b>92.08±0.14</b>	55.34±0.77	<b>86.12±0.99</b>
(BL-ZSL+SP-ZSR) [58]	69.74±3.47	92.06±0.18	53.26±1.04	86.01±1.32
<b>Ours: JFA+SP-ZSR [58]</b>	<b>80.89±5.97</b>	88.04±0.69	<b>55.81±1.37</b>	85.35±1.56

Table 4. Average precision and recall comparison (%) for recognition.

Precision	aP&Y	AwA	CUB	SUN
SDL [57]	52.70±27.33	81.70±14.67	54.06±24.13	82.51±12.24
<b>Ours: JFA</b>	52.41±25.59	79.31±11.70	53.73±23.90	85.12±10.87
SDL [57]+SP-ZSR [58]	55.96±35.72	<b>91.37±14.75</b>	57.09±27.91	85.96±10.15
<b>Ours: JFA+SP-ZSR [58]</b>	<b>76.41±29.70</b>	89.19±15.09	<b>57.20±25.96</b>	<b>86.06±12.36</b>
Recall				
SDL [57]	51.34±29.69	72.14±26.29	45.05±26.16	82.00±16.31
<b>Ours: JFA</b>	51.48±31.61	79.63±12.34	46.98±29.81	84.00±15.13
SDL [57]+SP-ZSR [58]	54.66±42.27	<b>90.28±8.08</b>	55.73±31.80	<b>86.00±13.19</b>
<b>Ours: JFA+SP-ZSR [58]</b>	<b>77.20±30.45</b>	86.04±9.82	<b>55.77±26.54</b>	85.50±13.68

lems with the help of attributes, because some bird species look visually very similar but still have their own unique characteristics. This leads to a more descriptive attribute vector per image than the average. In this context our adaptive matching tries to estimate the individual attribute vector from the average attribute vector of the class for matching based on the visual information. As we see on CUB, our method improves [57] by 4.70% in terms of accuracy, equivalently 11.25% relative improvement.

### 3.2.2 Transductive Setting

For transductive setting, we list our comparison results in Table 3. Overall, our method outperforms the state-of-the-art [58] significantly by 2.26% on average. It is worth mentioning that on aP&Y by substituting our similarity scores in [58] we can achieve 80.89% in terms of accuracy and outperform the state-of-the-art significantly by 11.15%. Analogous to the results of the traditional setting, we observe that the standard deviations of our results are slightly higher

than those of the competitors.

To better compare our results, we further measure the class-level performance on the datasets in terms of precision and recall (equivalent to accuracy per class). The detailed comparisons are illustrated in Fig. 5 without the CUB dataset due to the space limit. We summarize the averaged performance across different classes on each dataset in Table 4. Overall at the class level our method behaves similar to [57] with the same inputs. However, as we see there exists no single dominant method over all the datasets and uniformly over all classes on each dataset. Better similarity measure does not necessarily lead to better performance under either standard or transductive setting. It could be interesting as future work to see whether we can improve the ZSR performance further by integrating different similarity metrics.

## 4. Conclusion

In this paper we solve the relative sparseness issue of source-domain attribute vectors in ZSR problems. We formulate ZSL as a latent structural SVMs. To account for the rich data variability in target domain, we propose a novel data-dependent adaptive similarity function that adapts to test-time source and target data instances. Our similarity function searches for latent features from both domains by maximizing the latent similarities as well as minimizing the penalties incurred by feature displacements. To parameterize our adaptive similarity function, we propose a family of bilinear based similarity functions with regularized least squares to penalize displacements. We design a specific function with closed-form solutions and propose its corresponding learning algorithm for ZSR. To demonstrate the effectiveness of our proposed method, we test it on four benchmark datasets for ZSR with comprehensive comparison, and show significant improvement over the state-of-the-art under both standard and transductive settings.

## References

- [1] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid. Label-embedding for attribute-based classification. In *CVPR*, pages 819–826, 2013. [1](#), [2](#), [4](#), [5](#), [6](#), [7](#)
- [2] Z. Akata, S. Reed, D. Walter, H. Lee, and B. Schiele. Evaluation of output embeddings for fine-grained image classification. In *CVPR*, June 2015. [2](#), [4](#), [5](#), [6](#), [7](#)
- [3] Z. Al-Halah, M. Tapaswi, and R. Stiefelhagen. Recovering the missing link: Predicting class-attribute associations for unsupervised zero-shot learning. In *CVPR*, June 2016. [2](#)
- [4] S. Antol, C. L. Zitnick, and D. Parikh. Zero-shot learning via visual abstraction. In *ECCV*, pages 401–416. Springer, 2014. [1](#)
- [5] J. L. Ba, K. Swersky, S. Fidler, and R. Salakhutdinov. Predicting deep zero-shot convolutional neural networks using textual descriptions. *arXiv preprint arXiv:1506.00511*, 2015. [2](#)
- [6] T. L. Berg, A. C. Berg, and J. Shih. Automatic attribute discovery and characterization from noisy web data. In *ECCV*, pages 663–676, 2010. [1](#)
- [7] D. P. Bertsekas, A. E. Ozdaglar, and A. Nedi. *Convex analysis and optimization*. Athena scientific optimization and computation series. Athena Scientific, Belmont (Mass.), 2003. [4](#)
- [8] K. Bhatia, H. Jain, P. Kar, M. Varma, and P. Jain. Sparse local embeddings for extreme multi-label classification. In *NIPS*, 2015. [1](#)
- [9] M. Bucher, S. Herbin, and F. Jurie. Improving semantic embedding consistency by metric learning for zero-shot classification. *arXiv preprint arXiv:1607.08085*, 2016. [2](#), [7](#)
- [10] X. Chang, Y. Yang, A. G. Hauptmann, E. P. Xing, and Y.-L. Yu. Semantic concept discovery for large-scale zero-shot event detection. In *AAAI*, pages 2234–2240, 2015. [2](#)
- [11] S. Changpinyo, W.-L. Chao, B. Gong, and F. Sha. Synthesized classifiers for zero-shot learning. In *CVPR*, June 2016. [1](#), [2](#), [7](#)
- [12] 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. *arXiv preprint arXiv:1605.04253*, 2016. [2](#)
- [13] M. Elhoseiny, J. Liu, H. Cheng, H. Sawhney, and A. Elgammal. Zero-shot event detection by multimodal distributional semantic embedding of videos. *arXiv preprint arXiv:1512.00818*, 2015. [2](#)
- [14] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth. Describing objects by their attributes. In *CVPR*, pages 1778–1785, 2009. [1](#), [6](#), [7](#)
- [15] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *TPAMI*, 32(9):1627–1645, 2010. [4](#)
- [16] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. A. Ranzato, and T. Mikolov. Devise: A deep visual-semantic embedding model. In *NIPS*, pages 2121–2129, 2013. [1](#), [2](#)
- [17] Y. Fu, T. M. Hospedales, T. Xiang, Z. Fu, and S. Gong. Transductive multi-view embedding for zero-shot recognition and annotation. In *ECCV*, 2014. [3](#), [8](#)
- [18] Y. Fu, T. M. Hospedales, T. Xiang, and S. Gong. Transductive multi-view zero-shot learning. *PAMI*, 37(11):2332–2345, 2015. [3](#), [8](#)
- [19] Z. Fu, T. Xiang, E. Kodirov, and S. Gong. Zero-shot object recognition by semantic manifold distance. In *CVPR*, pages 2635–2644, 2015. [2](#)
- [20] C. Gan, M. Lin, Y. Yang, Y. Zhuang, and A. G. Hauptmann. Exploring semantic inter-class relationships (sir) for zero-shot action recognition. In *AAAI*, pages 3769–3775, 2015. [2](#)
- [21] Y. Guo, G. Ding, X. Jin, and J. Wang. Transductive zero-shot recognition via shared model space learning. In *AAAI*, 2016. [3](#), [8](#)
- [22] B. Hariharan, S. Vishwanathan, and M. Varma. Efficient max-margin multi-label classification with applications to zero-shot learning. *Machine learning*, 88(1-2):127–155, 2012. [2](#)
- [23] D. Jayaraman and K. Grauman. Zero-shot recognition with unreliable attributes. In *NIPS*, pages 3464–3472, 2014. [2](#), [6](#), [7](#)
- [24] T. Joachims, T. Finley, and C.-N. J. Yu. Cutting-plane training of structural svms. *Machine Learning*, 77(1):27–59, 2009. [6](#)
- [25] E. Kodirov, T. Xiang, Z. Fu, and S. Gong. Unsupervised domain adaptation for zero-shot learning. In *ICCV*, 2015. [2](#), [3](#), [8](#)
- [26] A. Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Master’s thesis, 2009. [6](#)
- [27] C. H. Lampert, H. Nickisch, and S. Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *CVPR*, pages 951–958, 2009. [1](#)
- [28] C. H. Lampert, H. Nickisch, and S. Harmeling. Attribute-based classification for zero-shot visual object categorization. *PAMI*, 36(3):453–465, 2014. [1](#), [2](#), [7](#)
- [29] X. Li and Y. Guo. Max-margin zero-shot learning for multi-class classification. In *AISTATS*, 2015. [2](#)
- [30] X. Li, Y. Guo, and D. Schuurmans. Semi-supervised zero-shot classification with label representation learning. In *ICCV*, 2015. [2](#)
- [31] D. Mahajan, S. Sellamanickam, and V. Nair. A joint learning framework for attribute models and object descriptions. In *ICCV*, pages 1227–1234, 2011. [2](#), [7](#)
- [32] T. Mensink, E. Gavves, and C. G. M. Snoek. Costa: Co-occurrence statistics for zero-shot classification. In *CVPR*, pages 2441–2448, June 2014. [2](#), [7](#)
- [33] 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*, pages 488–501, 2012. [1](#)
- [34] 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. [2](#)
- [35] M. Palatucci, D. Pomerleau, G. E. Hinton, and T. M. Mitchell. Zero-shot learning with semantic output codes. In *NIPS*, pages 1410–1418, 2009. [2](#)

- [36] D. Parikh and K. Grauman. Interactively building a discriminative vocabulary of nameable attributes. In *CVPR*, pages 1681–1688, 2011. 1
- [37] G. Patterson, C. Xu, H. Su, and J. Hays. The sun attribute database: Beyond categories for deeper scene understanding. *IJCV*, 108(1-2):59–81, 2014. 6
- [38] W. Ping, Q. Liu, and A. Ihler. Marginal structured svm with hidden variables. In *ICML*, pages 190–198, 2014. 3
- [39] R. Qiao, L. Liu, C. Shen, and A. van den Hengel. Less is more: Zero-shot learning from online textual documents with noise suppression. In *CVPR*, June 2016. 2, 7
- [40] M. Rohrbach, S. Ebert, and B. Schiele. Transfer learning in a transductive setting. In *NIPS*, pages 46–54, 2013. 2, 7
- [41] M. Rohrbach, M. Stark, and B. Schiele. Evaluating knowledge transfer and zero-shot learning in a large-scale setting. In *CVPR*, pages 1641–1648, 2011. 1
- [42] B. Romera-Paredes and P. H. S. Torr. An embarrassingly simple approach to zero-shot learning. In *ICML*, 2015. 2, 7, 8
- [43] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge, 2014. 1
- [44] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 6
- [45] R. Socher, M. Ganjoo, C. D. Manning, and A. Ng. Zero-shot learning through cross-modal transfer. In *NIPS*, pages 935–943, 2013. 1, 2
- [46] A. Vedaldi and K. Lenc. Matconvnet – convolutional neural networks for MATLAB. *CoRR*, abs/1412.4564, 2014. 6
- [47] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical report, 2011. 6
- [48] D. Wang, Y. Li, Y. Lin, and Y. Zhuang. Relational knowledge transfer for zero-shot learning. In *AAAI*, 2016. 2, 7
- [49] X. Wang and Q. Ji. A unified probabilistic approach modeling relationships between attributes and objects. In *ICCV*, pages 2120–2127, 2013. 2, 7
- [50] S. Wu, S. Bondugula, F. Luisier, X. Zhuang, and P. Natarajan. Zero-shot event detection using multi-modal fusion of weakly supervised concepts. In *CVPR*, pages 2665–2672, 2014. 1, 2
- [51] Y. Xian, Z. Akata, G. Sharma, Q. Nguyen, M. Hein, and B. Schiele. Latent embeddings for zero-shot classification. In *CVPR*, 2016. 2, 7
- [52] C.-N. J. Yu and T. Joachims. Learning structural svms with latent variables. In *ICML*, pages 1169–1176. ACM, 2009. 4
- [53] F. X. Yu, L. Cao, R. S. Feris, J. R. Smith, and S. F. Chang. Designing category-level attributes for discriminative visual recognition. In *CVPR*, pages 771–778, 2013. 1, 2, 7
- [54] X. Yu and Y. Aloimonos. Attribute-based transfer learning for object categorization with zero/one training example. In *ECCV*, pages 127–140. 2010. 2
- [55] Y. Zhang, B. Gong, and M. Shah. Fast zero-shot image tagging. In *CVPR*, June 2016. 2
- [56] Z. Zhang and V. Saligrama. Zero-shot learning via semantic similarity embedding. In *ICCV*, 2015. 1, 2, 6, 7
- [57] Z. Zhang and V. Saligrama. Zero-shot learning via joint latent similarity embedding. In *CVPR*, pages 6034–6042, 2016. 1, 2, 3, 4, 6, 7, 8
- [58] Z. Zhang and V. Saligrama. Zero-shot recognition via structured prediction. In *ECCV*, 2016. 3, 6, 8