

Attribute Embedding with Visual-Semantic Ambiguity Removal for Zero-shot Learning

Yang Long¹

ylong2@sheffield.ac.uk

Li Liu²

li2.liu@northumbria.ac.uk

Ling Shao²

ling.shao@ieee.org

¹ Department of Electronic and Electrical Engineering

The University of Sheffield
Sheffield, UK

² Department of Computer and Information Sciences
Northumbria University
Newcastle upon Tyne, UK

Abstract

Conventional *zero-shot learning* (ZSL) methods recognise an unseen instance by projecting its visual features to a semantic space that is shared by both seen and unseen categories. However, we observe that such a one-way paradigm suffers from the *visual-semantic ambiguity* problem. Namely, the semantic concepts (e.g. attributes) cannot explicitly correspond to visual patterns, and vice versa. Such a problem can lead to a huge variance in the visual features for each attribute. In this paper, we investigate how to remove such semantic ambiguity based on the observed visual appearances. In particular, we propose (1) a novel latent attribute space to mitigate the gap between visual appearances and semantic expressions; (2) a dual-graph regularised embedding algorithm called *Visual-Semantic Ambiguity Removal* (VSAR) that can simultaneously extract the shared components between visual and semantic information and mutually align the data distribution based on the intrinsic local structures of both spaces; (3) a new zero-shot recognition framework that can deal with both instance-level and category-level ZSL tasks. We validate our method on two popular zero-shot learning datasets, AWA and aPY. Extensive experiments demonstrate that our proposed approach significantly performs the state-of-the-art methods.

1 Introduction

Zero-shot learning focuses on classification with no training data. The fundamental idea of ZSL is to train a closed-set of human knowledge models that can generalise to an ever growing set of classes without collecting new training data. Such a scenario effectively alleviates the cost of data collection and also provides a feasible solution for recognising inaccessible objects, such as an ancient species that only has text records. Because of these attractive properties, ZSL has aroused increasing research interests in the vision and learning community [0, 21, 23].

Conventional ZSL methods [4, 8, 15, 25, 26] rely on directly mapping the visual features to a human-interpretable semantic space and the labels are inferred through human knowledge. However, an inevitable issue of using semantic information is the *ambiguity* problem. In linguistics, a concept is considered ambiguous if its extension is deemed lacking in clarity. It is the uncertainty about which objects belong to the concept or which exhibit characteristics that have this predicate. In the context of ZSL, **Visual-Semantic Ambiguity** refers to the situation that a semantic concept (e.g. an attribute) cannot clearly correspond to a certain pattern of visual data, and vice versa. Therefore, the paradox is how different of the visual patterns can we tolerate for each semantic concept? Alternatively, should we split the concept into sub-concepts to fit the visual data? This is known as the Sorites Paradox that can lead to two extreme solutions. (1) We can accept all instances as if they have the same attribute. Jayaraman and Grauman [11] also study this problem. They provide an extreme example that the concept ‘bumpy’ is assigned to both ‘bumpy road’ and ‘bumpy rash’ which can lead to unreasonable classification results. Unfortunately, most of the existing methods accept this solution. (2) We refuse any ambiguity and give every seen instance a unique attribute. For example, compared to ‘smile’, ‘Mona Lisa’s smile’ is clearly referring to a unique visual pattern with no ambiguity. However, it is infeasible to treat everything as unique and assign a new concept to it.

Instead of debating on what is common or unique, in this paper, we propose a latent attribute space to mitigate the visual-semantic ambiguity using a novel algorithm named *Visual-Semantic Ambiguity Removal* (VSAR). We measure the visual-semantic ambiguity by the reconstruction error and correct it in the latent attribute space. Intuitively, if a semantic concept refers to multiple variations of visual features, it should be split into different regions in the latent attribute space. In the visual aspect, if two close feature points are labelled by different attributes, we should find lower-dimensional subspaces so that they can be discriminated after embedding. Specifically, we develop a graph regularised embedding function that can minimise the reconstruction errors in both visual and semantic spaces. Meanwhile, the regularisation can preserve the discriminative information for recognising unseen categories. We illustrate this idea in Fig. 1. Our contribution is three-fold: (1) we propose a novel VSAR algorithm that can simultaneously remove the ambiguity between visual and semantic information; (2) our results suggest the important role of visual-semantic ambiguity to the performance improvement; (3) we introduce a unified framework that can deal with both category-level (AwA dataset) and instance-level (aPY dataset) zero-shot recognition tasks without adjusting the paradigm. To the best of our knowledge, the visual-semantic ambiguity issue has not been well studied yet. Thus, in the following, we only review related ZSL approaches.

Related work. Since learning visual attributes [5] is proposed, extensive studies [12, 13, 21, 25] have been conducted on how to use attributes as an intermediate representation for ZSL tasks. One interesting direction is to investigate the properties of attributes, such as the label co-occurrence property [20], the relativity [22], the unreliability [10], and the correlation problem [11] of human-nameable attributes. All of these are semantic properties and therefore suffer from the semantic-visual ambiguity problem. Due to this problem, some work turns to abandon human-nameable attributes and discovers data-driven attributes [14, 26]. However, for ZSL, these methods cannot exploit existing attribute ontologies. Hence, the applicable area is limited. Another trend is based on the embedding framework [1, 2, 3, 7, 19, 23, 28]. All these methods follow the restricted one-way paradigm that suffers

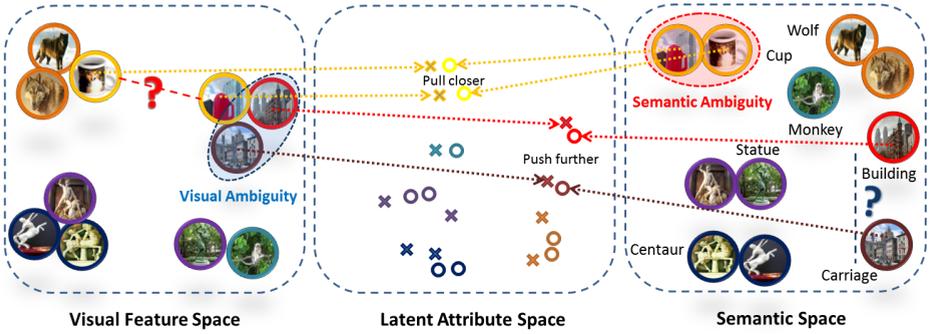


Figure 1: An intuitive illustration of VSAR (best viewed in colour). Visual Ambiguity (in blue oval): the image of a carriage is taken with a building background. It cannot recover the semantic distance (blue question mark) to the building category. Semantic Ambiguity (in red oval): the cup printed with a wolf and the cup-like building share the same semantic expression which can lead to a large visual variance (the red question mark). After embedding to the latent attribute space using VSAR, such ambiguity is mitigated.

from the ambiguity between low-level instances and high-level semantic concepts and labels. Recently, a new direction of ZSL is using the transductive model [6, 9, 13, 16, 18, 24, 27]. Unlabelled target domain data is collected for learning a transfer function. However, this setting slightly differs from the original ZSL purpose because the target domain may be strictly inaccessible. In contrast, our method can exploit the extensive existing attribute ontology while also stressing the existence of visual-semantic ambiguity and removing it through a learning process.

2 Visual-Semantic Ambiguity Removal

Problem setup: The training data is in N 3-tuples of ‘seen’ samples, attributes, and category labels: $(x_1, a_1, y_1), \dots, (x_N, a_N, y_N) \subseteq \mathcal{X}_s \times \mathcal{A}_s \times \mathcal{Y}_s$, where \mathcal{X}_s is a D -dimensional feature space $\mathcal{X}_s = [x_{dn}] \in \mathbb{R}^{D \times N}$, \mathcal{A}_s is a M -dimensional attribute space $\mathcal{A}_s = [a_{mn}] \in \mathbb{R}^{M \times N}$, and $y_n \in \{1, \dots, C\}$ consists of C discrete categories. The Calligraphic typeface indicates a space. We use subscript u to denote information of ‘unseen’ space and *hat* denotes ‘unseen’ samples. During testing, the preliminary knowledge is in \hat{C} pairs of ‘unseen’ category-level attributes and labels: $(\hat{a}_1, \hat{y}_1), \dots, (\hat{a}_{\hat{C}}, \hat{y}_{\hat{C}}) \subseteq \mathcal{A}_u \times \mathcal{Y}_u$, $\mathcal{Y} \cap \mathcal{Y}_u = \emptyset$, $\mathcal{A}_u = [a_{m\hat{c}}] \in \mathbb{R}^{M \times \hat{C}}$. The goal is to learn a classifier, $f: \mathcal{X}_u \rightarrow \mathcal{Y}_u$, where the samples in \mathcal{X}_u are completely unavailable during training. Such a problem is known as zero-shot learning.

Latent Attribute Embedding: We aim to discover a latent attribute embedding space \mathcal{V} shared by both visual and semantic spaces \mathcal{X} and \mathcal{A} to mitigate the visual-semantic ambiguity. During testing, both \mathcal{X}_u and \mathcal{A}_u can be embedded into \mathcal{V} .

Zero-shot Recognition: Instead of typical two-step prediction $\mathcal{X}_u \rightarrow \mathcal{A}_u \rightarrow \mathcal{Y}_u$, our embedding is two-way from \mathcal{X}_u and \mathcal{A}_u . Because attribute space \mathcal{A}_u and label space \mathcal{Y}_u are in pairs, we can firstly embed the known \mathcal{A}_u to \mathcal{V} as a knowledge domain. During testing, an unseen image \hat{x} is also embedded to \mathcal{V} so that we can compute the index, i.e., $\mathcal{X}_u \rightarrow \mathcal{V} \leftarrow \mathcal{A}_u \leftarrow \mathcal{Y}_u$.

2.1 Latent Attribute Embedding

This is the core component to deal with the visual-semantic ambiguity. We require \mathcal{X}_s and \mathcal{A}_s to compute \mathcal{V} . In the following, we drop the subscript s for convenience, i.e. we replace $\{\mathcal{X}_s, \mathcal{A}_s, \mathcal{Y}_s\}$ by $\{\mathcal{X}, \mathcal{A}, \mathcal{Y}\}$. Typically, each dimension a_m denotes a human-nameable concept, where $M \ll D$. The attribute notions here are instance-level. For the category-level, we can simply set the same attribute vectors to the instances within the same class. For embedding, many previous works are based on a forward matrix transformation, i.e. \mathcal{X} to \mathcal{A} . However, because of the visual-semantic ambiguity, the variance in \mathcal{X} is large. Therefore, the forward embedding is difficult to be reconstructed by a backward inverse matrix transformation from \mathcal{A} . Therefore, we insert an intermediate latent attribute space \mathcal{V} between \mathcal{X} and \mathcal{A} , where $\mathcal{V} = [v_{kn}] \in \mathbb{R}^{K \times N}$. K is the dimension of the embedding space. A straightforward setting is $M \leq K \leq A$. However, we stress that K can be any positive whole number. Specifically, we introduce our loss function as:

$$J = \|\mathcal{X} - U_1 \mathcal{V}\|_F^2 + \alpha \|\mathcal{A} - U_2 \mathcal{V}\|_F^2, \quad (1)$$

where $\|\cdot\|_F$ is the Frobenius norm of a matrix, which estimates the Euclidean distance between two matrices. The shared embedding space \mathcal{V} is decomposed from both \mathcal{X} and \mathcal{A} , where $U_1 = [u_{1dk}] \in \mathbb{R}^{D \times K}$ and $U_2 = [u_{2mk}] \in \mathbb{R}^{M \times K}$ are the basis matrices of the visual feature and attribute space, respectively.

Using Eq. 1, it becomes easier to understand the properties of the latent attribute space and how it could mitigate the visual-semantic ambiguity. Optimising Eq. 1 aims to minimise the reconstruction errors that are from \mathcal{V} to \mathcal{X} and from \mathcal{V} to \mathcal{A} , respectively. To achieve the optimal solution, U_1 and U_2 should preserve the principal components between \mathcal{X} and \mathcal{A} . This differs from unsupervised methods, such as PCA, that only analyse the data structure in a single domain. Our Eq. 1 can reduce the variance of the embedded data that comes from both visual and semantic domains. α is a reliability parameter that can balance the strengths of the two terms. In practice, if the attribute space is known as unreliable in prior, e.g. extended from category-level attributes, we can reduce α so that the proposed embedding can focus more on the visual feature space and remove more ambiguity from the attribute space.

2.2 Dual-graph Regularisation

The above Eq. 1 can reduce the difference between the data structures of \mathcal{X} and \mathcal{A} . However, it cannot preserve the discriminative information. For instance, if the gap between x_n and a_n is too large, their corresponding weights tend to be minimised to very small values. As a result, the learnt latent attributes are the principal components that are shared by all of the categories. For the purpose of zero-shot recognition, we have to preserve the intrinsic geometrical structure so that the learnt representation is discriminative.

We achieve this goal by taking the local invariance assumption and model the problem through a spectral graph approach named *Dual-graph Regularisation*. In particular, this is a combination of two supervised graphs that model the relationship between \mathcal{X} and \mathcal{Y} , and \mathcal{A} and \mathcal{Y} . The main criteria is to preserve the local structures. Therefore, we need the two graphs to simultaneously estimate the data structures of both spaces. Each graph has N vertices that correspond to N data points in the training set. As mentioned earlier, our method

can effectively handle ZSL tasks for both instance-level and category-level attribute scenarios. In particular, for *instance-level attributes*, we put an edge between each data point x_n or a_n and its p nearest neighbours. For each pair of the vertices s_i and s_j in the weight matrix, $w_{ij} = 1$ if and only if s_i and s_j are connected by an edge, otherwise, $w_{ij} = 0$. As a result, we can separately compute two weight matrices $W_{\mathcal{X}}$ and $W_{\mathcal{A}}$.

It is noteworthy that for *category-level attributes*, $W_{\mathcal{A}}$ is computed slightly different. Every vertex in the same category are connected by a normalised edge, i.e. $w_{ij} = p/n_c$, if and only if a_i and a_j are from the same category c , where n_c is the size of category c .

In the embedding space \mathcal{V} , we expect that if the s_i and s_j in both graphs are connected, each pair of embedded points v_i and v_j are also closed to each other. However, for the *visual-semantic ambiguity* problem, $W_{\mathcal{X}}$ and $W_{\mathcal{A}}$ usually give contradictory results. To compromise such conflict, we use the same reliability parameter α in Eq. 1 to linearly combine the two graphs, i.e. $W_{ij} = W_{\mathcal{X}_{ij}} + \alpha W_{\mathcal{A}_{ij}}$. The resulted regularisation is:

$$\begin{aligned} \mathcal{R} &= \frac{1}{2} \sum_{i,j=1}^N \|v_i - v_j\|^2 w_{ij} \\ &= Tr(\mathcal{V}D\mathcal{V}^T) - Tr(\mathcal{V}W\mathcal{V}^T) = Tr(\mathcal{V}L\mathcal{V}^T), \end{aligned} \quad (2)$$

where D is the degree matrix of W , $D_{ii} = \sum_i w_{ij}$. L is known as graph Laplacian matrix $L = D - W$ and $Tr(\cdot)$ computes the trace of a matrix. We combine Eq. 1 and 2 using a regularisation parameter λ to control the balance between reconstruction error and local structure preservation. The final goal is to optimise the following equation:

$$J = \|\mathcal{X} - U_1\mathcal{V}\|_F^2 + \alpha\|\mathcal{A} - U_2\mathcal{V}\|_F^2 + \lambda Tr(\mathcal{V}L\mathcal{V}^T), \quad (3)$$

2.3 Optimisation Strategy

Each term of the above Eq. 3 is convex, but the combined expression of U_1, U_2, \mathcal{V} is non-convex. To our best knowledge, there is no direct solution to find the global optima. Instead, we adopt an alternating optimisation strategy to find the local minima for each term separately as a relaxed solution. Specifically, the whole task is in turn separated into three sub-problems.

1. sub-problem U_1 : Suppose we compute the partial derivative of the overall loss function J with respect to U_1, U_2 and \mathcal{V} are fixed as constants. It then becomes a standard least squares problem. Let the partial derivative equal to zero, we have the closed form solution:

$$\begin{aligned} \frac{\partial J}{\partial U_1} &= -2\mathcal{X}\mathcal{V}^T + 2U_1\mathcal{V}\mathcal{V}^T = 0 \\ U_1 &= \mathcal{X}\mathcal{V}^T (\mathcal{V}\mathcal{V}^T)^{-1}. \end{aligned} \quad (4)$$

2. sub-problem U_2 : Similar to the sub-problem 1, we can fix U_1 and \mathcal{V} , and compute the partial derivative of J with respect to U_2 . The corresponding solution is:

$$U_2 = \mathcal{A}\mathcal{V}^T (\mathcal{V}\mathcal{V}^T)^{-1}. \quad (5)$$

Since we do not expect any prior bias from the unnormalised magnitudes of the training data, the basis vectors in the matrices should be normalised to unit vectors via:

$$u_{1_{dk}} \leftarrow \frac{u_{1_{dk}}}{\sqrt{\sum_d u_{1_{dk}}^2}} \quad u_{2_{mk}} \leftarrow \frac{u_{1_{mk}}}{\sqrt{\sum_m u_{2_{mk}}^2}}.$$

3. sub-problem \mathcal{V} : Fix U_1 and U_2 , we can then update \mathcal{V} . Applying the matrix properties $Tr(AB) = Tr(BA)$ and $Tr(A^T) = Tr(A)$, and we set the partial derivative respect to \mathcal{V} to zero:

$$\frac{\partial J}{\partial \mathcal{V}} = 2((U_1^T U_1 + \alpha U_2^T U_2) \mathcal{V} + \mathcal{V}(\lambda L) - (U_1^T \mathcal{X} + \alpha U_2^T \mathcal{A})) = 0. \quad (6)$$

Since space U_1 , U_2 and L are disjointed, this forms a typical Sylvester equation that has the unique solution for \mathcal{V} . We use the `lyap()` function in MATLAB to solve this problem.

Batch sampling scheme: In practice, the computational complexity of solving the Eq. 6 is $\mathcal{O}(N^3)$. To improve the efficiency, we adopt a batch sampling scheme like the deep learning strategy. The whole training set is divided into t batches by randomly sampling training instances from each categories. The size of each batch roughly equals to $\frac{N}{t}$. As a result, the computational complexity is reduced to $\mathcal{O}\left(t\left(\frac{N}{t}\right)^3\right)$, where $\left(\frac{N}{t}\right)^3 \ll N^3$. Each batch is in turn used to optimise the loss function in Eq. 3. We turn to the next batch until it converges on the previous batch. The whole learning procedure is summarised in Algorithm 1.

Algorithm 1 : Visual-Semantic Ambiguity Removal

Input: $\{\mathcal{X}, \mathcal{A}, \mathcal{Y}\}$, α, λ, K, p , number of batch t .

Output: \mathcal{V}, U_1 , and U_2 .

- 1: Initialisation: random batch sampling $\{\mathcal{X}_1, \mathcal{A}_1, \mathcal{Y}_1\} \dots \{\mathcal{X}_i, \mathcal{A}_i, \mathcal{Y}_i\}$, random initial matrix \mathcal{V} .
 - 2: **for** each batch **do**
 - 3: Compute the graph Laplacian matrix L using Eq. 2;
 - 4: **while** Eq. 3 is not converged **do**
 - 5: Update U_1 by Eq. 4, then normalise U_1 by $u_{1_{dk}} \leftarrow \frac{u_{1_{dk}}}{\sqrt{\sum_d u_{1_{dk}}^2}}$;
 - 6: Update U_2 by Eq. 5, then normalise U_2 by $u_{2_{mk}} \leftarrow \frac{u_{1_{mk}}}{\sqrt{\sum_m u_{2_{mk}}^2}}$;
 - 7: Update \mathcal{V} by Eq. 6;
 - 8: **end while**
 - 9: **end for**
 - 10: **return** \mathcal{V}, U_1 , and U_2 ;
-

2.4 Zero-shot Recognition

Once we obtain the latent attribute embedding \mathcal{V} of the seen data, performing zero-shot recognition is straightforward via *least-square approximation* between \mathcal{V} and $\{\mathcal{A}, \mathcal{X}\}$. During the test, the given informations are the unseen category names and their attributes in pairs: $\{\mathcal{Y}_u, \mathcal{A}_u\}$. We firstly embed all unseen attributes \mathcal{A}_u into the latent embedding space as references: $\mathcal{V}_u = \mathcal{V} \mathcal{A}^T (\mathcal{A} \mathcal{A}^T)^{-1} \mathcal{A}_u$. Given a test unseen instance \hat{x} , its embedded latent attribute representation is: $\hat{v} = \mathcal{V} \mathcal{X}^T (\mathcal{X} \mathcal{X}^T)^{-1} \hat{x}$. Finally, we adopt a simple NN classifier to

| Method | aPascal&Yahoo | Animals with Attributes |
|------------------------------|---------------------|-------------------------|
| Farhadi <i>et al.</i> [9] | 32.5 | - |
| Mahajan <i>et al.</i> [19] | 37.93 | - |
| Akata <i>et al.</i> [10] | - | 43.5 |
| Fu <i>et al.</i> [11] | - | 47.1 |
| Lampert <i>et al.</i> [15] | 19.1 | 40.5 |
| Jayaraman and Grauman [16] | 26.02 ± 0.05 | 43.01 ± 0.07 |
| Romera-Paredes and Torr [23] | 27.27 ± 1.62 | 49.30 ± 0.21 |
| our VSAR | 39.42 ± 0.27 | 51.75 ± 0.43 |

Table 1: Compare with the published state-of-the-art methods.

predict the category label \hat{c} :

$$\hat{c} = \arg \min_c \|\hat{v} - v_c\|^2, \text{ where } v_c \in \mathcal{V}_u. \quad (7)$$

3 Experiments

Datasets and Settings. We choose two of the most popular datasets for evaluating ZSL tasks. (a) **AwA dataset** [19] is one of the earliest work that particularly proposed for ZSL tasks. Many published results are based on this dataset. Each animal category in AwA is labelled by an attribute signature. (b) **aPY dataset** [9] is an instance-level attribute dataset that each image has a unique attribute signature. In contrast to AwA, aPY covers a more various range of categories, including human, artificial objects, buildings, as well as animals. For comparison reason, we adopt the base features that are provided by the datasets. We carefully follow the standard settings on both of the datasets. In particular, the training/test splits are 40/10 and 20/12 on AwA and aPY dataset, respectively. The optimal reliability parameter α for each dataset is selected from one of $\{0.1, \dots, 0.5, \dots, 0.9\}$ with the step of 0.1 which yields the best performance by 10-fold cross-validation on the training data. For λ and p , cross-validation is still deployed and finally fixed as $\lambda = 0.03$ and $p = 10$.

3.1 Comparison with the state-of-the-arts

We summarise our comparison in Table 1, where the hyphen indicates the existing method has not tested on the datasets in their original publication. Our method significantly outperforms the previous published results and can achieve state-of-the-art performance comparing to most recent papers. From the confusion matrices in Fig. 2 we can see that the recognition rate to each category tends to be averaged. Such a result indicates the performance of our proposed method is stable and reliable. It is also worth noting that, due to the attributes of the two datasets are not both category-level or instance-level, all of the compared methods have to adjust the framework to fit such different settings. In comparison, our VSAR approach can deal with both of the situations.

3.2 Algorithm analysis

Effects of terms in VSAR. To understand the success of our VSAR algorithm, the first important question is how does each terms in our VSAR algorithm work for ZSL. Thus, we separately strip-down each term in Eq. 3 into three baseline models. The first model is referred as *X-to-A*, in which we remove the second term of Eq. 3 and let the visual space \mathcal{X} directly map to the semantic space, i.e. $\mathcal{V} = \mathcal{A}$. This is exactly a DAP procedure that,

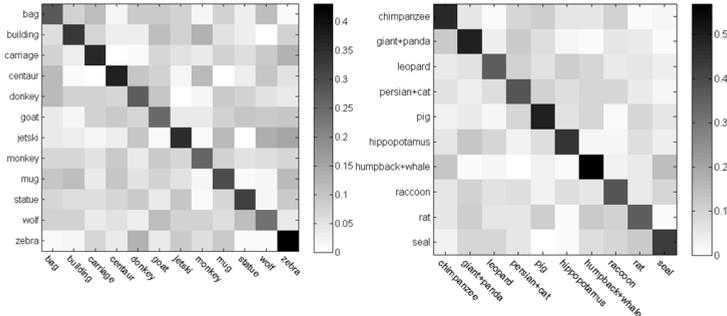


Figure 2: Confusion matrix of ZSL performance on aPY (left) and AWA (right).

during the test, the image is firstly mapped to the semantic space and then classified to the label space. The second model is referred as *A-to-X*. This is an interesting scenario that investigates whether we could regenerate the original visual features given just the semantic representations. Specifically, we train the model by setting $\mathcal{V} = \mathcal{X}$ and remove the first term in Eq. 3. During the test, we firstly project all attributes of the unseen categories/instances to \mathcal{X} . A test image is then classified in this embedding space using Eq. 7. In the third model that is denoted as *No-Graph*, we explore the importance of our dual-graph regularisations. Specifically, we train the model by setting $\lambda = 0$.

In Fig. 3, it can be seen that our full model significantly outperforms all of the baseline methods. In addition, we find the performance of the third model is roughly equal to random guess. Such a failure case matches our previous expectation that, without regularisation, Eq. 1 tend to discover the principle components rather than discriminating the categories. It is also noticeable that the *A-to-X* method gets better result on the aPY dataset than that on AWA. We ascribe this to the instance-level attributes. Such a result implies that it is feasible to generate visual features of each image from its semantic representations in future work.

Number of latent attributes. Another important issue is how many latent attributes K are required for the embedding space. Does a larger number of K always give better results? To investigate this question, we gradually increase K from 50, 85, 500, 1000, and 1000 per further step. We show the result in Fig. 3 (left). Generally speaking, a larger K tends to benefit the performance. However, we point out that there is an optimal K that gives the peak result. After that, the performance gradually degrades while we further increase K . This problem is severer on AWA than that on aPY. This is because when K goes too large, this can be viewed as an spectral over-fitting problem [29]. Since the attributes of AWA is category-level, the variance of its semantic space is much smaller than its visual space, which results in that the model on the AWA is more likely to over-fitting. For the whole experiments, we fix $K = 3000$ and 4000 for aPY and AWA, respectively.

Efficiency Our implementation is conducted in Matlab 2014a environment that is installed on a 12-core Linux system with 400G memory. The test time is done within a second. The training process takes roughly half an hour (i.e. number of batches $t = 15$) to get a converged model. Most of the time is used for solving the Eq. 6. We stress our contribution of using the batch sampling scheme, whereas directly solving the Eq. 7 without the batch sampling scheme can take up to 10 hours.

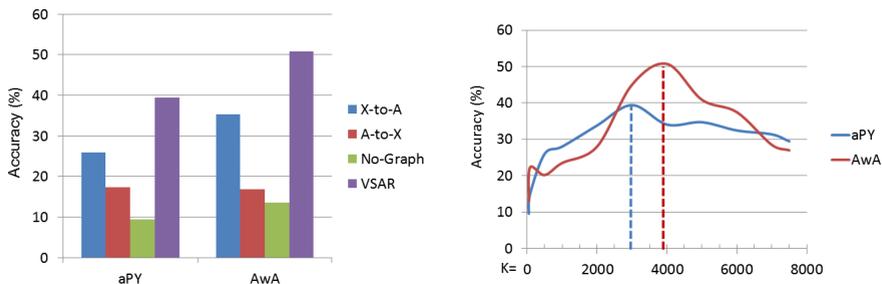


Figure 3: Evaluating each term of the loss function in Eq. 3 (left) and the performance curve respects to the dimension K of the latent attribute space (right).

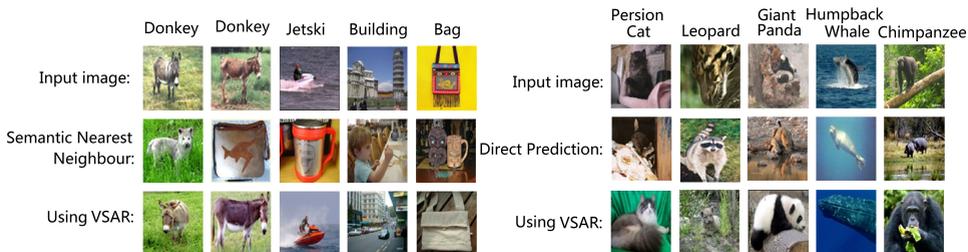


Figure 4: Examples of successful semantic ambiguity removal on aPY (left) and the visual ambiguity removal on AwA (right).

3.3 Visual-semantic ambiguity removal

In this section, we investigate what kinds of visual-semantic ambiguity are removed using our algorithm. This question can be considered from two aspects. Firstly, we consider the semantic ambiguity between different categories. On aPY dataset, we find such a semantic ambiguity problem is very severe. We use the provided “ground truth” attribute labels as the representation for each image. We then search the nearest neighbour for each image like an 1-NN classification. We find that only 67.17% of the nearest neighbours can match their original categories. Such a result implies that even if the conventional attribute classifiers can give perfect predictions, the overall recognition rate is only 67.17%. In Fig. 4 (left), we show that our VSAR is able to remove some of the semantic ambiguities. For example, in the second columns, the test image ‘donkey’ is misclassified as a ‘bag’ because the material and the logo of the bag possesses the same attributes to the donkey. However, in the visual space, such two instances are very distinctive. Therefore, using VSAR, our method successfully removes the ambiguity and gives the correct nearest neighbour. On the AwA dataset, the semantic ambiguity does not exist because all of the images in one category share the same attributes. Therefore, we consider the problem of visual ambiguity, i.e. the extracted low-level features from different categories are confused to each other. Specifically, we compare our method with the DAP framework using the X-to-A model. In Fig. 4 (right), we show some prediction errors in DAP can be corrected using VSAR. Such an ability contributes to the remarkable performance improvement (39.42% to 51.75%) in Fig. 3.

4 Conclusion and future work

We introduce that the visual-semantic ambiguity is a common issue in ZSL tasks. Our results on both datasets support that ambiguity removal can significantly benefit the recognition performance. The proposed VSAR is a unified framework that can deal with various semantic inputs, such as category-level and instance-level attributes. Instead of treating ZSL as a multi-label classification task, we adopt an embedding approach without struggling with the effectiveness of each attribute concept. Due to this property, our method can be simply applied to various existing intermediate semantic representations, such as data-driven attribute [26] or word-vector [25]. In the future, we plan to extend our visual-semantic constraints to multilateral in order to simultaneously incorporate multiple types of visual, semantic, as well as hierarchical label information.

References

- [1] Zeynep Akata, Florent Perronnin, Zaid Harchaoui, and Cordelia Schmid. Label-embedding for attribute-based classification. In *CVPR*, 2013.
- [2] Ziad Al-Halah, Tobias Gehrig, and Rainer Stiefelhagen. Learning semantic attributes via a common latent space. In *VISAPP*, 2014.
- [3] Ziyun Cai, Li Liu, Mengyang Yu, and Ling Shao. Latent structure preserving hashing. In *BMVC*, 2015.
- [4] Alireza Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. In *CVPR*, 2009.
- [5] Vittorio Ferrari and Andrew Zisserman. Learning visual attributes. In *NIPS*, 2007.
- [6] Yanwei Fu, Timothy M Hospedales, Tao Xiang, Zhenyong Fu, and Shaogang Gong. Transductive multi-view embedding for zero-shot recognition and annotation. In *ECCV*. 2014.
- [7] Yanwei Fu, Timothy M Hospedales, Tao Xiang, and Shaogang Gong. Learning multi-modal latent attributes. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 36(2):303–316, 2014.
- [8] Zhenyong Fu, Tao Xiang, Elyor Kodirov, and Shaogang Gong. Zero-shot object recognition by semantic manifold distance. In *CVPR*, 2015.
- [9] Yuchen Guo, Guiguang Ding, Xiaoming Jin, and Jianmin Wang. Transductive zero-shot recognition via shared model space learning. In *AAAI*, 2016.
- [10] Dinesh Jayaraman and Kristen Grauman. Zero-shot recognition with unreliable attributes. In *NIPS*, 2014.
- [11] Dinesh Jayaraman, Fei Sha, and Kristen Grauman. Decorrelating semantic visual attributes by resisting the urge to share. In *CVPR*, 2014.
- [12] Pichai Kankuekul, Aram Kawewong, Sirinart Tangruamsub, and Osamu Hasegawa. Online incremental attribute-based zero-shot learning. In *CVPR*, 2012.

- [13] Elyor Kodirov, Tao Xiang, Zhenyong Fu, and Shaogang Gong. Unsupervised domain adaptation for zero-shot learning. In *ICCV*, 2015.
- [14] Neeraj Kumar, Alexander C Berg, Peter N Belhumeur, and Shree K Nayar. Attribute and simile classifiers for face verification. In *ICCV*, 2009.
- [15] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *CVPR*, 2009.
- [16] Xin Li, Yuhong Guo, and Dale Schuurmans. Semi-supervised zero-shot classification with label representation learning. In *ICCV*, 2015.
- [17] Li Liu, Mengyang Yu, and Ling Shao. Projection bank: From high-dimensional data to medium-length binary codes. In *ICCV*, 2015.
- [18] Yang Long, Fan Zhu, and Ling Shao. Recognising occluded multi-view actions using local nearest neighbour embedding. *Computer Vision and Image Understanding*, 144: 36–45, 2016.
- [19] Dhruv Mahajan, Sundararajan Sellamanickam, and Vinod Nair. A joint learning framework for attribute models and object descriptions. In *ICCV*, 2011.
- [20] Thomas Mensink, Efstratios Gavves, and Cees Snoek. Costa: Co-occurrence statistics for zero-shot classification. In *CVPR*, 2014.
- [21] Mark Palatucci, Dean Pomerleau, Geoffrey E Hinton, and Tom M Mitchell. Zero-shot learning with semantic output codes. In *NIPS*, 2009.
- [22] Devi Parikh and Kristen Grauman. Relative attributes. In *ICCV*, 2011.
- [23] Bernardino Romera-Paredes and PHS Torr. An embarrassingly simple approach to zero-shot learning. In *ICML*, 2015.
- [24] Ling Shao, Li Liu, and Mengyang Yu. Kernelized multiview projection for robust action recognition. *International Journal of Computer Vision*, 118(2):115–129, 2016.
- [25] Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. Zero-shot learning through cross-modal transfer. In *NIPS*, 2013.
- [26] Felix Yu, Liangliang Cao, Rogerio Feris, John Smith, and Shih-Fu Chang. Designing category-level attributes for discriminative visual recognition. In *CVPR*, 2013.
- [27] Mengyang Yu, Li Liu, and Ling Shao. Structure-preserving binary representations for rgb-d action recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 38(8):1651–1664, 2016.
- [28] Ziming Zhang and Venkatesh Saligrama. Zero-shot learning via semantic similarity embedding. In *ICCV*, 2015.
- [29] Zheng Zhao and Huan Liu. Spectral feature selection for supervised and unsupervised learning. In *ICML*, 2007.