



Sharif University of Technology
Computer Engineering Department
MSc Thesis

Deep Zero-Shot Learning

Seyed Mohsen Shojaee

supervised by
Dr.Mahdieh Soleymani

Summer 2016

Agenda

1 Introduction

- Standard Learning Paradigm
- Zero-shot Learning definition
- Solution Steps

2 Prior Works

- Attribute Prediction
- Mapping to image space

3 Proposed Methods

- Multi-task Neural Network
- Mapping to Histogram of Seen Classes
- joining picture and lists
- pictures which need more space

Plan

1 Introduction

- Standard Learning Paradigm
- Zero-shot Learning definition
- Solution Steps

2 Prior Works

- Attribute Prediction
- Mapping to image space

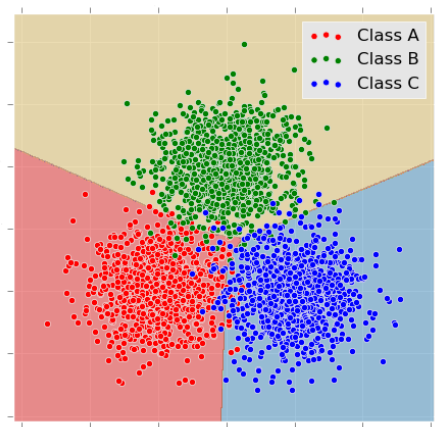
3 Proposed Methods

- Multi-task Neural Network
- Mapping to Histogram of Seen Classes
- joining picture and lists
- pictures which need more space

Introduction

Standard Learning Paradigm: Discover the pattern for each class from abundant labeled samples.

- Using SVM, Decision Tree, KNN, etc.

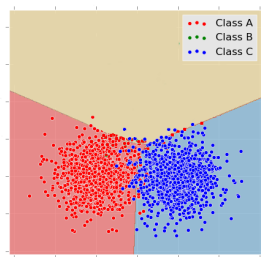


Extending the Standard Paradigm

Sometimes samples from all classes is not available

- Example: Novel Categories, Fine-grained classification.

Zero-Shot Learning addresses the problem of classification.
when no training sample is available for some classes.



Extending the Standard Paradigm

Identifying Classes without Samples:

- Each category is identified some *auxiliary information* also called *signature*.
- Examples of class signatures include:
 - Attribute Vectors
 - Text Articles
 - Category Names

Extending the Standard Paradigm

As a sample, an animal species like Zebra can have these signatures:

- The Vector (four legs, fast, striped, gallops, non-domestic, ...).
- The Wikipedia Entry for zebra.
- The word '*Zebra*' itself.



Problem Definition

At training time:

- there are N_s labeled samples: $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{N_s}$.
- These samples are from n_s classes that are called *seen classes*.
- Class signatures C_s for seen classes is also available.
- There are also n_u classes with no labeled sample. These are called *unseen classes*.
- It is assumed in most works that signatures of unseen classes, C_u , is also available.

Problem Definition

At test time:

- N_u samples from unseen classes are presented: $\{(\mathbf{x}_i)\}_{i=N_s+1}^{N_s+N_u}$.
- The Goal is to classify test samples into unseen categories.
- In other words finding

$$\arg \min_{\mathbf{y}^*_i} \mathbf{y}^*_i \neq \mathbf{y}_i, \quad i = N_s + 1, \dots, N_s + N_u$$

Solution Steps

Most existing solutions for Zero-shot learning consist of these three steps:

- 1 Embed images in a semantic space
- 2 Embed class signatures to same semantic space

Solution Steps

Most existing solutions for Zero-shot learning consist of these three steps:

- ① Embed images in a semantic space
- ② Embed class signatures to same semantic space
- ③ Assign images to those classes (e.g. using nearest neighbor classifier)

Solution Steps

Most existing solutions for Zero-shot learning consist of these three steps:

- ① Embed images in a semantic space
- ② Embed class signatures to same semantic space
- ③ Assign images to those classes (e.g. using nearest neighbor classifier)

Plan

1 Introduction

- Standard Learning Paradigm
- Zero-shot Learning definition
- Solution Steps

2 Prior Works

- Attribute Prediction
- Mapping to image space

3 Proposed Methods

- Multi-task Neural Network
- Mapping to Histogram of Seen Classes
- joining picture and lists
- pictures which need more space

Prior Works

Existing works can be categorized by the semantic space they use:

- Space of signatures (Attribute Prediction).
- Space of images.
- A third space.

We review some selected works from each category.

Attribute Prediction

- A large body of work in Zero-shot learning belongs to this category.
- The mapping from signature space is considered identity mapping.
- Attribute Estimator/Classifier are learned on train images (standard supervised problem).
- The Estimator/Classifier is used on test images to find \mathbf{c}_i^* for image \mathbf{x}_i
- \mathbf{x}_i is assigned to class with most similar signature:

$$\ell(\mathbf{x}_i) = \arg \min_{j=n_s+1, \dots, n_s+n_u} distance(\mathbf{c}_i^*, \mathbf{c}_j)$$

Mapping to Image Space

- In training time, Learn a mapping from class signatures to image space:

$$\phi : \mathbb{R}^a \rightarrow \mathbb{R}^d$$

- This can be seen predicting linear one-vs-all classifier for each class from its signature.
- In test time, classify test images using classifiers predicted from unseen class signatures.
- Assign each sample to class whose classifier produces maximum score:

$$\ell(\mathbf{x}) = \arg \max_{j=n_s+1, \dots, n_s+n_u} \langle \phi(\mathbf{c}_j), \mathbf{x} \rangle \quad (1)$$

Semi-supervised Zero-shot Learning

Plan

1 Introduction

- Standard Learning Paradigm
- Zero-shot Learning definition
- Solution Steps

2 Prior Works

- Attribute Prediction
- Mapping to image space

3 Proposed Methods

- Multi-task Neural Network
- Mapping to Histogram of Seen Classes
- joining picture and lists
- pictures which need more space

Proposed Methods

Here we present four proposed methods for the problem of Zero-shot Image Classification.

In our methods we consider class signatures of type attribute vectors.

- Attribute Prediction with Multi-task Deep Neural Networks.
- Mapping to Histograms of Seen Classes with Deep Neural Network.
- Embedding and Clustering
- Joint Embedding and Clustering

Multi-task Neural Network

We propose a network architecture for attribute prediction from images.

The network:

- predicts for train and test images at the same time (hence multi-task).
- can mitigate the domain shift problem that appears when only samples from seen classes is used.
- uses 16 convolutional layers from famous VGG-19 network [Simonyan and Zisserman, 2014] for feature extraction.
- is trained fast using Stochastic gradient descent algorithms family.

Multi-task Neural Network

Let f denote the mapping modeled by the multi-task network.

Then $\hat{\mathbf{c}}_i = f(\mathbf{x}_i)$ would be attributes predicted by network for \mathbf{x}_i

We learn f such that:

$$\underset{f}{\text{minimize}} \quad \frac{1}{N_s} \sum_{i=1}^{N_s} \text{loss}(\hat{\mathbf{c}}_i, \mathbf{c}_{y_i}) + \frac{\gamma}{N_u} \sum_{i=N_s+1}^{N_s+N_u} \left(\min_{j=n_s, \dots, n_s+n_u} \|\hat{\mathbf{c}}_i - \mathbf{c}_j\|_2^2 \right). \quad (2)$$

The second term enforces that prediction for test samples to be close to an unseen class signature

Therefore, mitigating domain-shift problem

Multi-task Neural Network

The second term in Eq. (2) is modeled by two layers, q and r :

$$(q(\mathbf{v}))_j = \|f(\mathbf{v}) - \mathbf{c}_j\|_2^2, \quad (3)$$

$$r(\mathbf{z}) = \min_{j=1 \dots n_u} (\mathbf{z})_j. \quad (4)$$

- The j -th element of q shows distance of prediction made by network to signature of j -th unseen category.
- r selects the minimum element of its input
- Hence using q and r successively produces distance of prediction to nearest unseen class signature.
- This is exactly same as the second term in Eq. (2)

Multi-task Neural Network

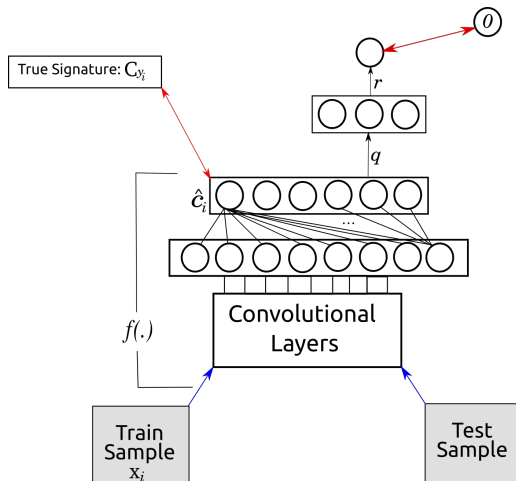


Figure: Proposed Multi-task network Architecture

Mapping to Histogram of Seen Classes

- Motivated by good performance of methods using histogram of similarity to seen classes as semantic space [Zhang and Saligrama, 2015].
- We present a deep neural network that maps images to this space.
- This network also uses convolutional layers from VGG-19 network.
- The network is a modification of a typical CNN used in standard supervised classification problems.

Mapping to Histogram of Seen Classes

- The network has a standard sequential architecture consisting of 17 pre-trained layers from VGG-19 and four other fully connected layers.
- Size of last layer is equal to the number of seen categories.
- Let ϕ denote the mapping modeled by the network

Mapping to Histogram of Seen Classes

In Training Time:

- Labeled samples from seen classes is used.
- Activation function in last layer is softmax:

$$\text{softmax}(\mathbf{z})_j = \frac{e^{z_j}}{\sum_k e^{z_k}}, \quad j = 1, \dots, n_s. \quad (5)$$

- Training criteria is correct label prediction of labeled samples.

$$\underset{\phi}{\text{minimize}} \sum_{i=1}^{N_s} \sum_{j=1}^{n_s} (\mathbf{y}^i)_j \times \log(\phi(\mathbf{x}_i)_j) + (1 - (\mathbf{y}^i)_j) \times \log(1 - \phi(\mathbf{x}_i)_j) \quad (6)$$

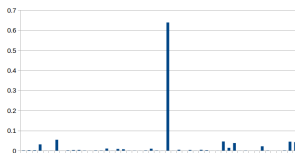
Mapping to Histogram of Seen Classes

In Test Time:

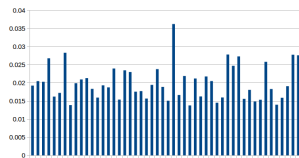
- Activation function in last layer is *temperature softmax*:

$$\text{softmax}_T(\mathbf{z})_j = \frac{e^{z_j/T}}{\sum_k e^{z_k/T}}, \quad T > 1, \quad j = 1, \dots, n_s. \quad (7)$$

- The softmax layer is trained to produce distribution of true label which is a discrete delta function.
- When setting $T > 1$ the output becomes smoother.



(a) $T = 1$



(b) $T = 10$

References I



Simonyan, K. and Zisserman, A. (2014).

Very deep convolutional networks for large-scale image recognition.

CoRR.



Zhang, Z. and Saligrama, V. (2015).

Zero-Shot Learning via Semantic Similarity Embedding.

In *IEEE Conference on International Computer Vision (ICCV)*, pages 4166–4174.