

A Shared Multi-Attention Framework for Multi-Label Zero-Shot Learning

Dat Huynh
Northeastern University
huynh.dat@husky.neu.edu

Ehsan Elhamifar
Northeastern University
eelhami@ccs.neu.edu

Abstract

In this work, we develop a shared multi-attention model for multi-label zero-shot learning. We argue that designing attention mechanism for recognizing multiple seen and unseen labels in an image is a non-trivial task as there is no training signal to localize unseen labels and an image only contains a few present labels that need attentions out of thousands of possible labels. Therefore, instead of generating attentions for unseen labels which have unknown behaviors and could focus on irrelevant regions due to the lack of any training sample, we let the unseen labels select among a set of shared attentions which are trained to be label-agnostic and to focus on only relevant/foreground regions through our novel loss. Finally, we learn a compatibility function to distinguish labels based on the selected attention. We further propose a novel loss function that consists of three components guiding the attention to focus on diverse and relevant image regions while utilizing all attention features. We show that our method improves the state of the art by 2.9% and 1.4% F1 score on the NUS-WIDE and the large scale Open Images datasets, respectively.

This is a short version of a full paper that is accepted for presentation in CVPR 2020.

1. Introduction

Multi-label learning is an important problem in image understanding since it captures a wide variety of objects appearing in the same image along with their correlation. It has been shown that using attention mechanism, which learns to localize labels under weakly supervision, significantly improves multi-label performance. In this work, we explore the possibility of generalizing attention mechanism to multiple unseen labels in an image. This setting is relevant because of the difficulty of collecting all possible samples during training. Therefore, the ability to extrapolate to novel classes not only significantly reduces annotation but also improve the system robustness when deploying in practice where novel classes could be encountered.

In this paper, we develop a framework for multi-label

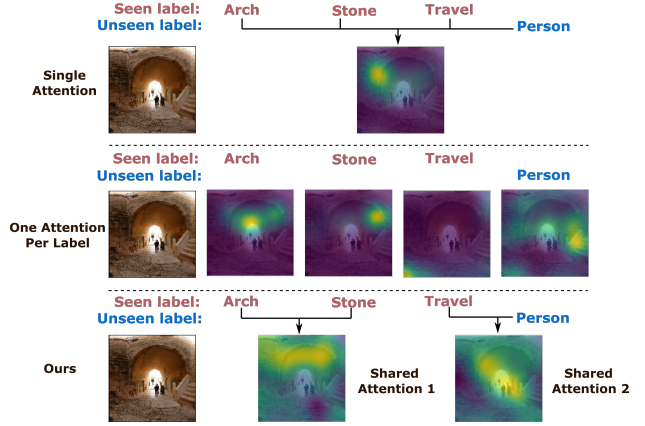


Figure 1: Visualization of attentions learned by single attention for all labels, one attention per label and our shared multi-attention model. Our shared multi-attention model successfully attends to relevant image regions for unseen labels while producing only a few number of attentions that significantly improves the memory and computational complexity for predicting thousands of labels.

zero-shot learning based on a novel shared multi-attention mechanism that handles recognition of a large number of labels, can recognize multiple unseen labels in an image and finds relevant regions to each label. Our method consists of multiple label-agnostic attention modules that generate multiple attention features simultaneously and uses the semantic vector of each label to select the most suitable feature to compute the prediction score of the label. Thus, instead of generating one attention feature for all labels, which cannot encode discriminative information about labels, and instead of generating one attention feature per label, which cannot generalize well to unseen labels, we generate multiple shared attention features that capture both common and discriminative information about labels, hence not only do well for the prediction of seen labels, but also transfer attention to unseen labels as in Figure 1.

Given that each training image only contains the list of present labels without ground-truth bounding box information, to effectively train our shared multi-attention method, we propose a novel loss function that enforces i) different attention modules to focus on diverse regions of an image,

covering different labels; ii) to find relevant regions that would lead to high prediction scores for present labels; iii) to effectively use all attention modules.

2. Multi-Label Zero-Shot Learning via Attention Sharing

Let $\mathcal{C}_s, \mathcal{C}_u$ be the set of seen and unseen classes respectively and $|\cdot|$ denote the cardinality of the set. We propose a shared multi-attention mechanism that consists of $M \ll |\mathcal{C}_s \cup \mathcal{C}_u|$ attention modules generating M attention features, where each feature will be used for the prediction of a subset of related labels, which are determined automatically. We also propose an efficient learning scheme that uses label semantic vectors and training images containing seen labels without access to their ground-truth localization. For an image i , where we divide the image into R equal regions and extract its region features $\{\mathbf{f}_i^r\}_{r=1}^R$, let $\{\mathbf{z}_i^m\}_{m=1}^M$ denote M attention features obtained via the attention modules $\{g_m(\cdot)\}_{m=1}^M$. We define

$$\mathbf{F}_i \triangleq [\mathbf{f}_i^1 \dots \mathbf{f}_i^R], \quad \boldsymbol{\alpha}^m(\mathbf{F}_i) \triangleq [\alpha_1^m(\mathbf{f}_i^1) \dots \alpha_R^m(\mathbf{f}_i^R)]^\top, \quad (1)$$

where \mathbf{F}_i denotes a matrix whose columns are R region features and $\boldsymbol{\alpha}^m(\mathbf{F}_i)$ denotes the R -dimensional weight vector of the attention module m , for the image i . we formulate the m -th attention feature of the image i , denoted by \mathbf{z}_i^m , as a linear combination of all region features as

$$\mathbf{z}_i^m = \mathbf{F}_i \boldsymbol{\alpha}^m(\mathbf{F}_i). \quad (2)$$

To learn and infer $\boldsymbol{\alpha}^m(\mathbf{F}_i)$, we use a simple two-layer neural network model

$$\boldsymbol{\alpha}^m(\mathbf{F}_i) = \frac{\exp(\mathbf{e}_i^m)}{\sum_{r=1}^R \exp(\mathbf{e}_{i,r}^m)}, \quad \mathbf{e}_i^m = \tanh(\mathbf{F}_i^\top \mathbf{W}_1^m) \mathbf{w}_2^m, \quad (3)$$

where $\{\mathbf{W}_1^m, \mathbf{w}_2^m\}_{m=1}^M$ are the model parameters, $\tanh(\cdot)$ is the element-wise hyperbolic tangent function, $\boldsymbol{\alpha}^m(\mathbf{F}_i)$ is the softmax normalization on each elements $\mathbf{e}_{i,r}^m$ of the R -dimensional unnormalized attention weights, \mathbf{e}_i^m , (i.e., before applying softmax) from the attention module m . Given M attention features $\{\mathbf{z}_i^m\}_{m=1}^M$, we propose a model in which the score of each label $c \in \mathcal{C}_s \cup \mathcal{C}_u$ is obtained by the maximum response of the classifier $\boldsymbol{\theta}^c$ (which is defined latter in Eq 10) over the M attention features, i.e.,

$$s_i^c \triangleq \max_{m=1, \dots, M} \langle \boldsymbol{\theta}^c, \mathbf{z}_i^m \rangle. \quad (4)$$

Thus, different attention features can be used for the prediction of different labels. To learn the parameters of the M attention modules, we propose an efficient learning scheme with a novel loss function, which we discuss next.

Diverse Multi-Attention Features: We ideally want different attention modules to attend to different regions of an

image. Thus, we define a diversity loss that promotes obtaining diverse attention features for an image. More specifically, using the cosine similarity between distinct pairs of unnormalized attention weight vectors, we define

$$\mathcal{L}_{div} \triangleq \sum_i \sum_{m \neq n} \frac{\langle \mathbf{e}_i^m, \mathbf{e}_i^n \rangle}{\|\mathbf{e}_i^m\|_2 \|\mathbf{e}_i^n\|_2}, \quad (5)$$

whose minimization promotes small or no overlap in the focus regions of different attention modules. For efficient learning, we use unnormalized attention weights \mathbf{e} instead of normalized weights $\boldsymbol{\alpha}$, since the gradient of $\boldsymbol{\alpha}$ vanishes when softmax function saturates. Also, we do not minimize the inner product $\langle \mathbf{e}_i^m, \mathbf{e}_i^n \rangle$, since it reduces not only the cosine similarity but also the ℓ_2 -norm of each weights vector, which prevents the weights of an attention module to concentrate on a single region.

Relevant Multi-Attention Features: Given that the training data does not include information about locations of labels in images, we cannot learn attention models by enforcing that attention weights on ground-truth regions be larger than weights on irrelevant regions. Here, we are only given the set of existing labels in each image. To tackle the problem, we use the prediction scores as surrogates for relevant regions to attend.

Our key observation is that when a seen label $o \in \mathcal{C}_s$ is present in an image, there must be a region containing o on which we have a high score for the label o . Thus, when successfully focused on the region of a label, the score of our multi-attention mechanism must be larger than simply weighting all regions equally. More specifically, let $\bar{s}_i^o \triangleq \frac{1}{R} \sum_r \langle \boldsymbol{\theta}^o, \mathbf{f}_i^r \rangle$ be the average score of the label o across all regions, i.e., the score when all regions contribute equally. We define a region relevance loss function that promotes our multi-attention mechanism produce higher scores than \bar{s}_i^o for present labels and lower scores for absent labels. In other words, we define

$$\mathcal{L}_{rel} \triangleq \sum_i \sum_{o \in \mathcal{C}_s} \max \left((\bar{s}_i^o - s_i^o) y_i^o, 0 \right), \quad (6)$$

where $y_i^o \triangleq 1$ for $o \in \mathcal{Y}_i$ and $y_i^o \triangleq -1$ otherwise. Notice that with the above loss, attention modules find not only regions of present labels, but also indicative regions of absent labels, e.g., to predict the absence of the label ‘desert’, the attention may focus on a region with the label ‘ocean’.

Using All Multi-Attention Modules: Given the ability to select among M different attention features in (4) and the non-convexity of learning, the model could potentially learn to use only some attention modules for prediction of all labels and not use the rest. Thus, we propose a loss function to encourage that each of the M attention modules will be used for the prediction of some of the seen labels. We start by defining a score ℓ_m that measures the utility of the m -th

attention module by computing the number of labels across training images that use the attention module m ,

$$\ell_m \triangleq \sum_i \sum_{o \in \mathcal{Y}_i} \mathbb{I}_m(\text{argmax}_n \langle \theta^o, z_i^n \rangle), \quad (7)$$

where $\mathbb{I}_m(x)$ is the indicator function, which outputs 1 when $x = m$ and 0 otherwise. Notice that the term inside the first sum in (7) corresponds to the number of labels of the image i that use the attention model m , hence, ℓ_m measures the utility of the attention module m across all training images. Ideally, we want every attention modules to be used for predictions, hence, we want to avoid having a few large ℓ_m 's while most being zero. Thus, we propose to minimize the attention distribution loss,

$$\mathcal{L}_{dist} \triangleq \sum_{m=1}^M \ell_m^2. \quad (8)$$

The difficulty of minimizing \mathcal{L}_{dist} is that the ℓ_m defined in (7) is non-differentiable, due to the indicator function. We tackle this by using a softmax function instead, where

$$\ell_m \triangleq \sum_i \sum_{o \in \mathcal{Y}_i} \frac{\exp(\langle \theta^o, z_i^m \rangle)}{\sum_{n=1}^M \exp(\langle \theta^o, z_i^n \rangle)}. \quad (9)$$

Notice that softmax function approximates the indicator of argmax, with the two coinciding when the magnitude of $\langle \theta^o, z_i^m \rangle$ is significantly larger than other $\langle \theta^o, z_i^n \rangle$.

Bilinear Compatibility Function: Given that we do not have training images for \mathcal{C}_u , we cannot directly optimize over and learn θ^u for $u \in \mathcal{C}_u$. Thus, we use the pre-trained GloVe representation of label names as semantic vectors $\{v^c\}_{c \in \mathcal{C}}$ of labels, allowing to transfer knowledge from seen to unseen labels. More specifically, we express the parameters of each classifier as a function of its semantic vector $\theta^c = \mathbf{W}_3 v^c$ and substituting in (4), compute the compatibility score of each label $c \in \mathcal{C}$ in an image i as

$$s_i^c = \max_{m=1, \dots, M} \langle \mathbf{W}_3 v^c, z_i^m \rangle. \quad (10)$$

Once we learn \mathbf{W}_3 , as discussed below, we can determine the labels in an image i by ranking and picking the top prediction scores $\{s_i^c\}_{c \in \mathcal{C}}$ across all labels.

To learn the parameters of the compatibility function, \mathbf{W}_3 , and the attention models, we use the ranking loss that imposes the scores of present labels in each image be larger by a margin than the scores of absent labels. More specifically, we define the ranking loss as

$$\mathcal{L}_{rank} \triangleq \sum_i \sum_{o \in \mathcal{Y}_i, o' \notin \mathcal{Y}_i} \max(1 + s_i^{o'} - s_i^o, 0), \quad (11)$$

in which the margin is set to one.

Final Loss Function: Putting all loss functions, discussed above, together we propose to minimize

$$\min_{\Theta, \{\mathbf{W}_1^m, \mathbf{w}_2^m\}_m, \mathbf{W}_3} \mathcal{L}_{rank} + \lambda_{div} \mathcal{L}_{div} + \lambda_{rel} \mathcal{L}_{rel} + \lambda_{dist} \mathcal{L}_{dist}, \quad (12)$$

where $\lambda_{div}, \lambda_{rel}, \lambda_{dist} \geq 0$ are regularization parameters. We minimize this loss function using stochastic gradient descent. In the experiments, we investigate the effectiveness of each loss function term.

3. Experiments

3.1. Experimental Setup

Datasets: We perform experiments on the NUS-WIDE [5] and the Open Images [6] datasets. In the NUS-WIDE, we use the 925 labels as seen and the other 81 labels as unseen as in [3] and follow the original training and testing splits in [5]. To demonstrate the effectiveness of our method on a larger number of labels and images, we use the large-scale Open Images (v4) dataset, which consists of 9 millions training images and 125,436 testing images. We use the 7186 training classes as the seen classes and select 400 disjoint classes from seen classes in test set as unseen classes.

Evaluation Metrics: We use the mean Average Precision (mAP) and F1 score at top K predictions in each image. Notice that the *mAP* score captures how accurate the model ranks images for each *label*, while the *F1* score measures how accurate the model ranks present labels in each *image*.

Implementation Details: Similar to other works [3], we use VGG-19 to extract the feature map of the last convolutional layer ($14 \times 14 \times 512$) as the region features. We implement all methods in Tensorflow and optimize with the default setting of RMSprop with the learning rate 0.001 and batch size of 32. We do not perform heavy hyperparameter tuning and set $(\lambda_{div}, \lambda_{rel}, \lambda_{dist})$ to $(1e^{-2}, 1e^{-3}, 1e^{-1})$ for both datasets. We set the number of attention modules to $M = 10$, unless stated otherwise. We refer to our method as L_Earning by Sharing Attentions (LESA).

3.2. Experimental Results

Multi-Label Zero-Shot Learning: We consider both multi-label zero-shot learning, where models are trained on seen labels and tested only on unseen labels, and multi-label generalized zero-shot learning, where models are tested on both seen and unseen labels. Table 1 shows the mAP score and F1 score at $K \in \{3, 5\}$ for NUS-WIDE and at $K \in \{10, 20\}$ for Open Images. We use a larger K for Open Images, since models need to make a larger number of predictions due to a much larger number of labels. From the results, we make the following observations:

- Our method outperforms the state of the art on both datasets, improving the mAP score on NUS-WIDE by 4.3%

