Interactive Multi-Label CNN Learning with Partial Labels

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

> 45600 68400 Number of iterations

Figure 1: Left: Improvement of mAP score (%) of our method with respect to Logistic regression for groups of classes with least annotations (left bar) to most annotations (right bar) in Open Images. The red curve shows the number of available training images in each group. Our method significantly improves over Logistic especially for the middle groups that have sufficient annotations. Right: mAP score of our method on the validation set during training. We observe that the performance curve fluctuates at the early stage of the training since the model struggles to refine noisy label graph and learn the data manifold. However, after 20k iterations, the learning curve monotonically improves until convergence. Here, each iteration corresponds to a gradient descent step on a minibatch.

1. Additional Results on Open Images

Figure 1 shows the mAP improvement of our method (IMCL) over the logistic regression when dividing labels into 10 groups (instead of 5, discussed in the main paper). Notice that our method improves the mAP score over the logistic regression for all the 10 groups. This shows the effectiveness of using the collaborative similarity learner in conjunction with the logistic model, which regularizes the network and takes advantage of semantic similarities across images and labels. Also, the right plot in Figure 1 shows the improvement of the mAP score on the validation set during training. We observe that the performance fluctuates at the early stage of the training as the model tries to refine the noisy label graph and learn the data manifold. However, after 20k iterations, the learning curve monotonically improves until convergence.

2. Detailed Results on CUB Dataset

Table 1 shows the exact performance of different methods as a function of the percentage of missing attributes on the CUB dataset (in the main paper, we showed bar charts). Ehsan Elhamifar Northeastern University

eelhami@ccs.neu.edu

Missing attributes	Logistic	CNN-RNN	Fast0Tag	Curriculum Labeling	$\begin{array}{l} \text{IMCL} \\ a = 1 \end{array}$	$\begin{array}{c} \text{IMCL} \\ a = 10 \end{array}$
90%	23.8	23.3	23.3	23.8	25.3	25.7
80%	25.2	25.0	25.0	25.7	26.4	26.4
60%	26.4	27.3	27.3	27.4	27.3	27.9
40%	26.6	26.9	27.2	27.8	27.7	27.8
20%	26.6	27.3	27.6	27.8	27.5	27.6
0%	27.9	27.2	27.9	27.9	28.6	28.7

Table 1: mAP scores (%) as a function of the percentage of missing attributes in the CUB dataset.

CNN-RNN	LSEP	Wsabie	Fast0Tag	Latent Noise (relevant)	Latent Noise (visual)	Curriculum Labeling	IMCL
-0.6	0.1	0.1	0.0	0.9	1.0	1.5	1.9

Table 2: Improvement of mAP score (%) of different methods over the logistic regression on the MS-COCO dataset.

3. Detailed Results on MS-COCO Dataset

Table 2 shows the exact mAP improvement of different methods over the logistic regression on the MS-COCO dataset (in the main paper, we showed bar charts).

4. Interactive Learning Framework

In this section, we present the derivations of our proposed Joint Nonnegative OMP algorithm for finding semantically similar images, in the similarity learning step in our framework, presented in the main paper. We then discuss the computational complexity of our algorithm.

4.1. Joint Nonnegative OMP Derivation

Consider our objective function with the classifier parameters, $\{\boldsymbol{\theta}_j\}_{j=1}^{|\mathcal{L}|}, w$, written as

$$\min_{\boldsymbol{w},\boldsymbol{\theta}_{1},\dots,\boldsymbol{\theta}_{|\mathcal{C}|}} \sum_{i} \mathcal{L}_{c}^{(i)} \Big(\boldsymbol{w}, \{\boldsymbol{\theta}_{j}\}_{j=1}^{|\mathcal{C}|} \Big) + \mathcal{L}_{s}^{(i)} \Big(\boldsymbol{w}, \{\boldsymbol{\theta}_{j}\}_{j=1}^{|\mathcal{C}|} \Big),$$
(1)

where the cross entropy loss for image i is defined as

$$\mathcal{L}_{c}^{(i)} \triangleq -\sum_{j \in \Omega_{i}} y_{j,i}^{o} \log(p_{j,i}) + (1 - y_{j,i}^{o}) \log(1 - p_{j,i}), \quad (2)$$

and the prediction smoothness loss for image i is defined as

$$\mathcal{L}_{s}^{(i)} \triangleq \min_{\{c_{i',i},\bar{c}_{i',i}\}} \lambda_{y} \left\| \boldsymbol{y}_{i} - \tanh\left(\sum_{i'=1}^{N} c_{i',i} \boldsymbol{A} \boldsymbol{y}_{i'}\right) \right\|_{2}^{2} \\ + \lambda_{f} \left\| \boldsymbol{f}_{i} - \sum_{i'=1}^{N} \bar{c}_{i',i} \boldsymbol{f}_{i'} \right\|_{2}^{2} \\ \text{s. t. } \sum_{j} \mathrm{I}\left(\left\| \left[c_{i',i}, \ \bar{c}_{i',i} \right] \right\| \right) \leq k, \ c_{i',i}, \bar{c}_{i',i} \geq 0, \ \forall i'. \end{cases}$$
(3)

To efficiently solve for image and label similarities, we use a first-order approximation of the hyperbolic tangent function in $\mathcal{L}_{s}^{(i)}$, which is $\tanh(x) \approx x$. This is a good approximation as long as \tanh is not saturated. Using this approximation, we can rewrite (3) as

$$\begin{aligned} \mathcal{L}_{s}^{(i)} &\triangleq \min_{\{c_{i',i}, \bar{c}_{i',i}\}} \lambda_{y} \| \boldsymbol{y}_{i} - \sum_{i'=1}^{N} c_{i',i} \boldsymbol{A} \boldsymbol{y}_{i'} \|_{2}^{2} \\ &+ \lambda_{f} \| \boldsymbol{f}_{i} - \sum_{i'=1}^{N} \bar{c}_{i',i} \boldsymbol{f}_{i'} \|_{2}^{2} \\ \text{s. t. } \sum_{j} \mathrm{I} \Big(\| \big[c_{i',i}, \ \bar{c}_{i',i} \big] \| \Big) \leq k, \ c_{i',i}, \bar{c}_{i',i} \geq 0, \ \forall i'. \end{aligned}$$

$$(4)$$

To compute the loss function $\mathcal{L}_{s}^{(i)}$, we need to solve the joint nonnegative sparse optimization in (4) over both $c_{i',i}$ and $\bar{c}_{i',i}$. We generalize [1], a greedy nonnegative sparse solver, to efficiently solve for both nonnegative similarity coefficients. More specifically, we use a generalization of OMP that starts from the empty set, at each iteration selects the best image i' that minimizes the smoothness loss the most (i.e., whose label and feature vectors best reconstruct the label and feature vectors of image i), updates residual errors and repeats to select the next best candidate, until k candidates are chosen. Let \mathcal{N}_i denote the set of similar images chosen so far. We compute the the optimal similarity coefficients { $c_{i',i}, \bar{c}_{i',i}$ } and the residual errors for the label, r_y , and feature $, r_f$, as

$$c_{i',i}^{*} = \max\left(0, \arg\min \|\boldsymbol{y}_{i} - \sum_{i' \in \mathcal{N}_{i}} c_{i',i} \boldsymbol{A} \boldsymbol{y}_{i'}\|_{2}^{2}\right),$$

$$\bar{c}_{i',i}^{*} = \max\left(0, \arg\min \|\boldsymbol{f}_{i} - \sum_{i' \in \mathcal{N}_{i}} \bar{c}_{i',i} \boldsymbol{f}_{i'}\|_{2}^{2}\right),$$

$$\boldsymbol{r}_{y} = \boldsymbol{y}_{i} - \sum_{i' \in \mathcal{N}_{i}} c_{i',i}^{*} \boldsymbol{A} \boldsymbol{y}_{i'},$$

$$\boldsymbol{r}_{f} = \boldsymbol{f}_{i} - \sum_{i' \in \mathcal{N}_{i}} \bar{c}_{i',i}^{*} \boldsymbol{f}_{i'}.$$
(5)

For each image i' not in the current set \mathcal{N}_i , we compute the loss $\mathcal{L}_s^{(i)}(\mathcal{N}_i \cup \{i'\})$ and then select the best i' for which

we have the minimum loss function. To do so, we fix the similarity coefficients for N_i and compute

$$\mathcal{L}_{s}^{(i)}(\mathcal{N}_{i} \cup \{i'\}) = \min_{c_{i',i}} \lambda_{y} \| \boldsymbol{r}_{y} - c_{i',i} \boldsymbol{A} \boldsymbol{y}_{i'} \|_{2}^{2} + \lambda_{f} \| \boldsymbol{r}_{f} - \bar{c}_{i',i} \boldsymbol{f}_{i'} \|_{2}^{2}$$
(6)

We select the image i' that achieves the minimum loss function, i.e., $s = \arg \min_{i'} \mathcal{L}_s^{(i)}(\mathcal{N}_i \cup \{i'\})$. For each i', we compute the closed-form of (6) by setting the derivative with respect to $c_{i',i}$ and $\bar{c}_{i',i}$ to zero,

$$\frac{\partial \mathcal{L}_{s}^{(i)}(\mathcal{N}_{i} \cup \{i'\})}{\partial c_{i',i}} = \frac{\partial \|\boldsymbol{r}_{y} - c_{i',i}\boldsymbol{A}\boldsymbol{y}_{i'}\|_{2}^{2}}{\partial c_{i',i}} = 0,$$

$$(\boldsymbol{A}\boldsymbol{y}_{i'})^{T} \Big(c_{i',i}^{*}\boldsymbol{A}\boldsymbol{y}_{i'} - \boldsymbol{r}_{y} \Big) = 0, \quad (7)$$

$$\implies c_{i',i}^{*} = \frac{\langle \boldsymbol{r}_{y}, \boldsymbol{A}\boldsymbol{y}_{i'} \rangle}{\|\boldsymbol{A}\boldsymbol{y}_{i'}\|_{2}^{2}}.$$

Similarly, we obtain the optimal value for $\bar{c}_{i',i}$ as

$$\bar{c}_{i',i}^* = \frac{\langle \boldsymbol{f}_{i'}, \boldsymbol{r}_f \rangle}{\|\boldsymbol{f}_{i'}\|_2^2}.$$
(8)

Substituting (7) and (8) into (6), we can compute the optimal loss function for any given $i', \mathcal{L}_{s}^{(i)}(\mathcal{N}_{i} \cup \{i'\})$ as

$$\begin{split} \lambda_{y} \| \boldsymbol{r}_{y} - \frac{\langle \boldsymbol{r}_{c}, \boldsymbol{A} \boldsymbol{y}_{i'} \rangle}{\| \boldsymbol{A} \boldsymbol{y}_{i'} \|_{2}^{2}} (\boldsymbol{A} \boldsymbol{y}_{i'}) \|_{2}^{2} + \lambda_{f} \| \boldsymbol{r}_{f} - \frac{\langle \boldsymbol{f}_{i'}, \boldsymbol{r}_{f} \rangle}{\| \boldsymbol{f}_{i'} \|_{2}^{2}} \boldsymbol{f}_{i'} \|_{2}^{2} \\ &= \lambda_{y} \Big(\frac{\langle \boldsymbol{r}_{y}, \boldsymbol{A} \boldsymbol{y}_{i'} \rangle^{2} \| \boldsymbol{A} \boldsymbol{y}_{i'} \|_{2}^{2}}{\| \boldsymbol{A} \boldsymbol{y}_{i'} \|_{2}^{2}} - 2 \frac{\langle \boldsymbol{r}_{y}, \boldsymbol{A} \boldsymbol{y}_{i'} \rangle^{2}}{\| \boldsymbol{A} \boldsymbol{y}_{i'} \|_{2}^{2}} + \| \boldsymbol{r}_{c} \|_{2}^{2} \Big) \\ &+ \lambda_{f} \Big(\frac{\langle \boldsymbol{r}_{f}, \boldsymbol{f}_{i'} \rangle^{2} \| \boldsymbol{f}_{i'} \|_{2}^{2}}{\| \boldsymbol{f}_{i'} \|_{2}^{2}} - 2 \frac{\langle \boldsymbol{r}_{f}, \boldsymbol{f}_{i'} \rangle^{2}}{\| \boldsymbol{f}_{i'} \|_{2}^{2}} + \| \boldsymbol{r}_{f} \|_{2}^{2} \Big) \\ &= -\lambda_{y} \frac{\langle \boldsymbol{r}_{y}, \boldsymbol{A} \boldsymbol{y}_{i'} \rangle^{2}}{\| \boldsymbol{A} \boldsymbol{y}_{i'} \|_{2}^{2}} - \lambda_{f} \frac{\langle \boldsymbol{r}_{f}, \boldsymbol{f}_{i'} \rangle^{2}}{\| \boldsymbol{f}_{i'} \|_{2}^{2}} + \text{constant.} \end{split}$$

Thus, we select the best next sample in Algorithm 2 via

$$s = \arg\max_{i'} \lambda_y \frac{\langle \boldsymbol{r}_y, \boldsymbol{A}\boldsymbol{y}_{i'} \rangle^2}{\|\boldsymbol{A}\boldsymbol{y}_{i'}\|_2^2} + \lambda_f \frac{\langle \boldsymbol{r}_f, \boldsymbol{f}_{i'} \rangle^2}{\|\boldsymbol{f}_{i'}\|_2^2}.$$
 (10)

4.2. Speeding up Training

The line 4 of Algorithm 2 in the main paper requires comparing the residual vectors with label vector and feature vector of every image. For large datasets, such as Open Images or MS-COCO, where there is a lot of redundancy in images, we could significantly reduce the computational time of this step by using a subset of images, obtained using random sampling or subset selection techniques. In our experiments, we construct the smaller dictionary by randomly selecting, for each label in the training set, 10 images that contain that label. Thus, the dictionary for Open Images and MS-COCO have the sizes of 50,000 and 10,000, respectively (we used the whole CUB dataset since there are only 11,000 images in the dataset). This simple strategy speeds up the training while improving the state-of-the-art results, as shown in the paper.

4.2.1 Complexity Analysis

For the Joint Nonnegative OMP in Algorithm 2 of the main paper, the dominant complexity cost of the algorithm is to find similarities (line 4) and solve the non-negative least squares (line 9 and 10). Let N_S be the size of the dictionary used to find similar images and let l and f be, respectively, the dimension of y_i and f_i . When the search for similar images is performed on the dictionary of size N_S , the complexity is $\mathcal{O}(N_S(l+f))$. On the other hand, solving non-negative least squares has the complexity of $\mathcal{O}(k^3 + k^2(l+f))$. Finding all k similar images thus takes $\mathcal{O}(N_S k(l+f) + k^4 + k^3(l+f))$. With $N_S \ll N$, as is in our experiments on the Open Images and MS-COCO datasets, the complexity of the Joint Nonnegative OMP for one image would be $\mathcal{O}(1)$ in the size of the dataset. Thus, finding similarities for all images has $\mathcal{O}(N)$ complexity. Note that (1) is optimized in a mini-batch fashion, where we iteratively construct images similarity for a small batch of data. Thus, we have low memory complexity, $\mathcal{O}(B)$ with B being the size of the minibatch, in terms of the size of the training data in each iteration.

References

 T. H. Lin and H. T. Kung, "Stable and efficient representation learning with nonnegativity constraints," *International Conference on Machine learning*, 2014. 2