

Informative Object Annotations - Supplementary Material

Lior Bracha

Bar-Ilan University

lior.bracha@live.biu.ac.il

Gal Chechik

Bar-Ilan University, NVIDIA Research

gal.chechik@biu.ac.il

Abstract

Supplementary materials to the main paper.

1. Robustness to hyper parameters

We tested the robustness of our approach with respect to two hyper parameters of the model: (1) number of trees on which we averaged to create Tree-Mixture-Model. (2) vocabulary size.

For the first, we computed all scoring functions for tree mixtures with 1,3,5 and 10 trees, and found only a 3% difference in the p@1 of $cw\Delta H$.

Second, we tested robustness to the number of words in our vocabulary. The vocabulary size is important because our analysis was performed over the most frequent labels in the corpus. As a result, the size of the vocabulary could have affected precision, because entry-level terms (*dog, car*) tend to be more frequent than more fine-grained terms (*e.g. Labrador, Toyota*). We repeated our analysis with different thresholds on the minimum label frequency included in the vocabulary (threshold for values of 50, 100:1000). Figure S 1 plots the precision@1 of the various scoring functions, showing that the analysis is robust to the size of the vocabulary.

2. Expected entropy

The *expected entropy-reduction* was defined as

$$E(\Delta H) = q(l_i|I)\Delta H + (1 - q(l_i|I)) \cdot 0$$

When an incorrect label is transmitted, we assume here that no information is passed to the listener. However, this does have an impact on the listener entropy and knowledge scheme. Transmission of a false message could be thought of as negative information as it is misleading for the listener.

3. Implementation Details

Algorithm 1 describes in detail the steps to compute the $cw\Delta H$ scores for a set of labels. Algorithm 2 describes

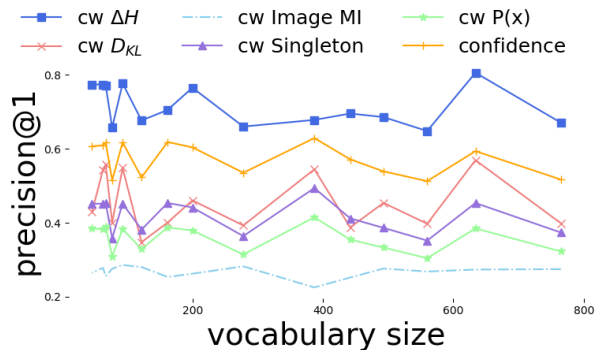


Figure S 1. **Robustness to vocabulary size.** Different thresholds for the minimum number of label occurrence were tested. The precision of $cw\Delta H$ remains very high for a large range of vocabulary sizes. The relationship between the different scoring functions is consistent as well.

the inference phase, where the computed scores provide an information-based ranking of the image annotations. Here, we do not specify whether we take a single label as ground truth (by majority) or multiple labels (see Sec 4.2) but give a general framework.

4. Qualitative Results

Figure S2 illustrates the annotations ranking for some images from OID test-set. In these examples we give the full, raw output of our experiments, showing results from all scoring-functions, with or without the confidence weights. "verification" column specifies whether the label was verified by OID raters as correct. "R*" columns present our raters response (see Sec. 4.1) and "y-true" column is the ground truth determined by majority. R* columns in which no entry is marked "1", means that the rater's label was not in the vocabulary.

Algorithm 1 $cw\Delta H$ scoring-function

Input Annotated images**Output** $cw\Delta H$ label-scores (averaged on k trees)

```
1: for tree do
2:   Sample uniformly from the label-distribution  $D$  (bootstrap)
3:   if  $label\ frequency > threshold$  then
4:      $vocab \leftarrow label$ 
5:   end if

6:   for label pair  $(i, j) \in vocab$  do
7:     Compute pair (2x2) joint distribution  $p(i, j)$ 
8:     Compute the mutual information of  $i, j$  (Eq. 7)
9:   end for

10:  Create Graph  $G(V, E)$ 
11:  Assign  $MI_{ij}$  as the weight of the edge connecting label  $i$  and label  $j$ 
12:  Find a maximum weight spanning tree (MST)
13:  Sort the graph such that each node has a single parent
14:  Compute tree entropy (Eq.8)

15:  for each label  $l_i$  in  $l_1, \dots, l_d$  do
16:    Set  $l_i$  as root
17:    Set root marginal distribution as  $p(l_i) = [0, 1]$ 
18:    Reverse edges direction such that all  $e \in E$  will be descendents of  $l_i$ 
19:    Propagate the message  $p(l_i)$  throughout the tree and update CPTs.
20:    Compute new tree entropy (Eq. 8)
21:    Compute  $\Delta H(l_i)$  (Eq. 3)
22:  end for
23: end for

24: Average  $\Delta H(l_i)$  over trees for all labels  $(l_1 \dots l_d)$ 
```

Algorithm 2 Ranking and evaluation

Input $\Delta H(l_i)$ label scores; ground-truth data**Output** Ranking of image annotations, precision and recall

```
1: for image do
2:    $cw\Delta H \leftarrow confidence(l_i) \cdot \Delta H(l_i)$ 
3:   Rank image annotations by  $cw\Delta H$ .
4:   Evaluate against ground-truth label.
5:   Compute precision and recall.
6: end for
7: Average precision and recall across images.
```
