

# Semantic-Analysis Object Recognition: Automatic Training Set Generation Using Textual Tags

Sami Abduljalil Abdulhak<sup>1</sup>(✉), Walter Riviera<sup>1</sup>, Nicola Zeni<sup>2</sup>,  
Matteo Cristani<sup>1</sup>, Roberta Ferrario<sup>2</sup>, and Marco Cristani<sup>1</sup>

<sup>1</sup> Department of Computer Science, Cá Vignal 2, Verona, Italy

{sami.naji,walter.riviera,matteo.cristani,marco.cristani}@univr.it

<sup>2</sup> Laboratory for Applied Ontology, Consiglio Nazionale Delle Ricerche (CNR),  
Via Alla Cascata 56/c, Trento, Italy

{nicola.zeni,roberta.ferrario}@loa.istc.cnr.it

**Abstract.** Training sets of images for object recognition are the pillars on which classifiers base their performances. We have built a framework to support the entire process of image and textual retrieval from search engines, which, giving an input keyword, performs a statistical and a semantic analysis and automatically builds a training set. We have focused our attention on textual information and we have explored, with several experiments, three different approaches to automatically discriminate between positive and negative images: keyword position, tag frequency and semantic analysis. We present the best results for each approach.

**Keywords:** Training set · Semantic · Ontology · Semantic similarity · Image retrieval · Textual tags · Flickr · Object recognition

## 1 Introduction

The process of automatically building a training set of images for object recognition given a class name is a recent challenge originated from the Semantic Robot Vision Challenge [1]. The idea is to mine on-line repositories of images and use them to support image classifiers in object recognition tasks [2]. Given this strategy, the goal is to exploit search engines and retrieve images that can be used to feed a training set for a specific class.

The problem falls under the topic of Image Retrieval (IR): given a certain query in a form of a keyword or an image, the system should present images related to the query. Two main strategies have been deployed to tackle such problem: content-based image retrieval (CBIR) [3] and tag/keyword-based image retrieval (TBIR)[4].

CBIR leverages on the concept of visual similarity between the querying image and the retrieved ones using elementary visual features such as color and shape, through a matching of their properties, while TBIR tries to overcome the

limitations presented by the CBIR system through the exploitation of the textual information conveyed with images, applying document retrieval techniques to boost the retrieval performances. Nevertheless TBIR performances are influenced by the availability and quality of the textual information users supply with images. In fact, while manually annotating images, users often misuse tags or provide incomplete textual descriptions of the image content [5–7].

The use of the textual information conveyed with images in the process of image retrieval or image classification is not a novel strategy, there have been several works that explore how the textual information can be used, among them [8–11]. Recent approaches explore the use of tags completion either by mining extra textual information obtained from Internet or by using content image analysis to fill the gap[6, 12].

In the present work we propose a framework that helps to automate the entire process of training data set construction. The main idea is to use textual information that comes along with images on the web to fully automate the training set generation. To achieve this, we assume that the user annotation process is not always reliable since users are not experts and may annotate images with different purposes. Even though users upload images in a social context where other users can use collaborative tagging to annotate images, tags are not validated and so the subjectivity elements are not removed. Moreover, since users are non expert, they tend to use ambiguous and inappropriate tags to describe images content. The main idea is to explore how statistical and semantic analysis of textual information can help to fully automate the training set construction. In particular, we employ statistical and semantic analysis to filter the textual information, pruning noisy tags and retaining only those that are highly correlated with the content of an image, thus discriminating positive from negative images<sup>1</sup>. We use statistical measures such as frequency and tags distribution, as well as WordNet and semantic distances between tags to evaluate their correlation and explore their contribution in the discriminative process. Our starting assumption is that, by incrementally injecting semantic techniques into the analysis of textual annotation, performances rise and, to validate such assumption, a set of experiments are presented.

The rest of the paper is structured as follows: Section 2 describes the challenges of the image retrieval task and provides an overview of works in the area. The method we propose is introduced in Section 3. Sections 4 discusses the experimental setup and evaluation method, while the evaluation results are presented in Section 5. Finally, conclusions and directions for future work are presented in Section 6.

---

<sup>1</sup> We consider as positive those images in which the prominence of the object presented in the image indicates that the image fully represents it. On the contrary, we consider as negative those images where the target object is absent or only partially present/visible, as indicated in the list in section 4.

## 2 Related Work

Annotation is a widely used technique to characterize objects portrayed in images by adding textual tags. The textual tags associated with images have been shown to be useful, improving the access to photo repositories both using temporal [13] and geographical information [14]. One of the popular online tag-based photo sharing repositories is Flickr, allowing users to freely assign one or more chosen keywords for an image for personal organization or retrieval purposes. In other words, it allows users to perform tagging, that is the act of adding words to images, describing the semantics of the visual contents. Users are thus implicitly encouraged to add more keywords, creating relatively large amounts of rich descriptions of objects presented in images. However, the textual tags associated with images are often noisy and unreliable, posing a number of difficulties when dealing with IR.

A number of approaches have been proposed to measure the reliability of the textual tags accompanying images [15–17]. In [17], the authors present a Flickr distance to measure the correlation between different concepts obtained from Flickr. Given a pair of concepts (e.g., car-dog), the algorithm tries to compute the semantic distance between them using square root of Jensen-Shannon divergence. The authors rely on the scores by considering the higher score distance as an indication of high relatedness of a pair of concepts. Related researches have been also focused on investigating which objects people observe most in an image, which they annotate or tag first, and what influence them in choosing words to describe objects depicted in images.

Spain and Perona [15] study the idea of “*importance*” of objects in an image and conclude that important objects are most likely to be tagged first by humans when asked to describe the contents of an image. The authors develop a statistical model validating the notion of dominant object in an image, demonstrating that one can foresee a set of prominent keywords based on the visual cues through regression. A work that is closely related to ours is presented by Hwang and Grauman [18]. They introduce an unsupervised learning method for IR that uncovers the implicit information about the object importance in an image, exploiting a list of keyword tags provided by humans. The proposed method is able to disclose the relationship between human tendencies in tagging images (e.g., words order in the tag list) and the relative importance of objects in an image.

Traditional techniques rely on features extracted from visual contents with visual category models learnt directly from image repositories that require no manual supervision [8–11]. The intuition behind the approach proposed in [9] is to learn object categories from just a few training images in an incremental manner, using a generative probabilistic model. Similarly, Li-Jia Li and Fei-Fei Li [10] propose an incremental learning framework, capable of automatically collecting large image datasets. The authors build a database from a sample of seed images and use the database to filter out newly crawling images by eliminating irrelevant examples.

Fergus et al. [11] introduce a method able to learn object categories by their name, exploiting the raw images automatically downloaded from the Google

image search engine. The introduced approach is able to incorporate spatial information in translation and scale invariant style, possessing the ability to tackle the high intra-category variability and isolate irrelevant images produced by the search engine.

Vijayanarasimhan and Grauman [19] propose an unsupervised approach to learn visual categories by their names using a collection of images pooled from keyword-based search engines. The main goal underneath the proposed approach is to harvest multiple images, by translating the query names into several languages and crawling the search engines for images using those translated queries. The false positive categories are collected from random sample images found in categories that have different names from the category of interest.

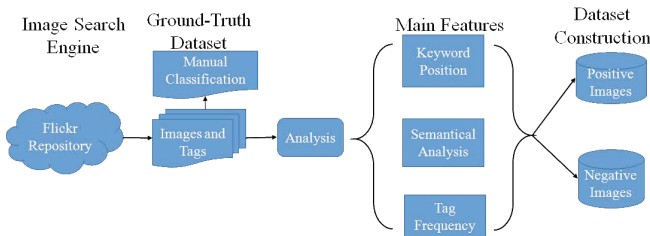
We are working on a challenge that is: given the textual tags provided by humans and associated with images, we want to automatically build a good training set by discriminating images as either related or unrelated to a targeted object.

### 3 Method

In this paper our goal is to take advantage of the textual tags available with images to automatically select the most representative of an object category for training a classifier, without looking at the nature of the objects therein. To do so, we exploit both semantic analysis and pure statistical approaches. These considerations lead us to focus on three main features:

- **keyword position**, to capture an image as related or unrelated on the basis of a keyword (i.e., object class name) position in a tag list;
- **semantic analysis**, to measure the semantic relatedness by means of semantic distance measures;
- **tag frequency**, to count the frequency of usage of each tag from a list describing the object class.

Figure 1 presents a schematic representation of our framework. A detailed description of the procedure is provided in the subsequent subsections.



**Fig. 1.** A schematic representation of our framework.

### 3.1 Keyword Position

The textual tags given in a tag list and associated with an image describing its content could reasonably help us to derive important and valuable information about the nature of the depicted objects. However, the order in which the textual tags are placed in a tag list is most likely to be influenced by the objects position and size in the visual content [20]. Therefore, it is reasonable to claim that the first textual tags in the list are mostly representing the objects in the center of an image. Taking this keypoint into account, we use this feature to develop 5 different strategies which follow the same algorithmic structure:

---

#### Algorithm 1. Keyword Position

---

**Data:** a Keyword (i.e., the object class name) and

$$T = \{t_i \mid \forall \text{ Image } i \in \text{Keyword}, \exists \text{ tag-list } t_i \}$$

**Result:** A partition of the Images  $\in$  Keyword in:

Image- $P = \{i_p \mid i \in \text{Images which are usable to build a training dataset}\}$

Image- $N = \{i_n \mid i \in \text{Images which are outliers}\}$

```

1 Initialization;
2 foreach  $i \in \text{Images}$  do
3    $\text{tags} \leftarrow \text{load } t_i$ ;
4    $\text{clean}(\text{tags})\text{tags}_n \leftarrow \text{extract the first } n \text{ tags from } \text{tags}$ ;
5   if "keyword"  $\in \text{tags}_n$  then
6     | Image- $P \leftarrow i$ ;
7   else
8     | Image- $N \leftarrow i$ ;

```

---

Algorithm 1 is designed to demonstrate the systematic workflow of the keyword position feature. Given a tag list comprising a number of textual tags and corresponding to a particular image, the algorithm tries to search for the keyword through the list in the first  $n$  positions. The algorithm then labels the image as positive (reliable) if it is related to the class name or negative (outlier) otherwise. It is noteworthy that the clean operation provided in the algorithm is used to remove words with less than three characters, empty strings and non-alphabetic texts. It also splits long sentences into single words, when they are separated by the “.” symbol.

### 3.2 Semantic Analysis

To define the semantic relatedness or its inverse of the object class characterized by a keyword to the textual tags being used, semantic distance must be measured. Therefore we propose to apply two different standard semantic distance measures: WordNet and Jiang and Conrath [20]. First we adopt the WordNet distance [21]. WordNet is a large-scale lexical database that organizes English terms and their syntactic roles into synsets. Synsets are interlinked by means of conceptual-semantic and a variety of lexical relations. We choose WordNet due to the fact that it is the first attempt to organize a great amount of concepts according to semantic relations and a hierarchy. Since WordNet provides

a lexical relationship between concepts, it is beneficial to semantically measure relatedness of the object class to its related tags by their lexical relationship, such as meronymy (parthood, e.g. bus-wheels) or hypernym (generalization, e.g. bus-vehicle) and so on.

Secondly, we apply the distance measure proposed by Jiang and Conrath in [20]. They formulate their approach in the form of conditional probability of reaching an item of a child synset given an item of one of its parent synsets.

We use this feature and run several experiments according to the following algorithmic structure:

---

**Algorithm 2.** Semantic Analysis

---

**Data:** a Keyword (i.e., the object class name) and

$$T = \{t_i \mid \forall \text{ Image } i \in \text{Keyword}, \exists \text{ tag-list } t_i \}$$

**Result:** A partition of the Images  $\in$  Keyword in:

Image- $P = \{i_p \mid i \in \text{Images which are usable to build a training dataset}\}$

Image- $N = \{i_n \mid i \in \text{Images which are outliers}\}$

1 Initialization;

2 **foreach**  $i \in \text{Images}$  **do**

3      $tags \leftarrow \text{load } t_i$ ;

4     **clean**( $tags$ );

5      $score_i \leftarrow$  sum or mean of the **distance** values of the  $tags$ ;

6     **if**  $if \text{ score}_i \geq a \text{ Threshold } \tau$  **then**

7         Image- $P \leftarrow i$ ;

8     **else**

9         Image- $N \leftarrow i$ ;

---

Algorithm 2 is developed to clearly illustrate how we apply the semantic analysis feature to measure the semantic relatedness or its inverse of the object class to its textual tags. As already mentioned above, we adopt two different distance measures: WordNet and Jiang and Conrath. The algorithm takes the object class (represented by a keyword) and each image's tag list, then computes the distance of the keyword to every single textual tag in the tag list, yielding a score for each. If the algorithm finds no semantic distance between the keyword and a textual tag, it discards the tag. The algorithm therefore labels an image as positive (reliable) if its score is equal or above a threshold  $\tau$ ; otherwise it labels it as negative (outlier). The threshold value  $\tau$  changes with respect to the experiment (see Section 4).

### 3.3 Tag Frequency

To understand which are the most frequently used tags (words) that describe images related to a certain object class, we compute the frequency values of all the single  $tag_{(i,j)}$  as their occurrences probability. The idea is to perform a selection based on the utility of the words used to describe the object depicted

in an image. The frequency value of a single  $tag_{(i,j)}$  is computed as follows:

$$Freq(tag_{(i,j)}) = \frac{O - tag_{(i,j)}}{\sum_{i=1}^{N_{images}} length(tag_i)},$$

where  $tag_{(i,j)}$  is the  $j^{th}$  tag of the tag list associated to image  $i$ , and  $O - tag_{(i,j)}$  is the total number of a  $tag_{(i,j)}$  occurrences. In particular, if a given frequency value of a single  $tag_{(i,j)}$  is relatively high, it means that many images of the considered object class require it into their descriptions. In other words, it is natural to think that if we are looking at an image of a “car”, we highly expect to observe higher frequency values for tags like “wheel” or “driver” than “pizza” or “pencil”.

We use this feature to develop 12 different strategies which follow the same algorithmic structure:

---

**Algorithm 3.** Tag Frequency

---

**Data:** a Keyword (i.e., the object class name) and

$$T = \{t_i | \forall Image\ i \in Keyword, \exists tag\ -\ list\ t_i\}$$

**Result:** A partition of the  $Images \in Keyword$  in:

Image- $P = \{i_p | i \in Images\ \text{which are usable to build a training dataset}\}$

Image- $N = \{i_n | i \in Images\ \text{which are outliers}\}$

```

1 Initialization;
2 foreach  $i \in Images$  do
3    $tags \leftarrow load\ t_i$ ;
4    $clean(tags)$ ;
5    $score_i \leftarrow$  sum or mean of the frequency values of the  $tags$ ;
6   if  $if\ score_i \geq a\ Threshold\ \tau$  then
7     | Image- $P \leftarrow i$ ;
8   else
9     | Image- $N \leftarrow i$ ;
```

---

Algorithm 3 uses frequency values to determine if a given image is related to the object class. To do this, it combines the frequency values of each  $tag_{(i,j)}$  to produce a score. Then, it labels an image  $i$  as positive (reliable) if its score is equal or above a threshold  $\tau$ ; otherwise it labels it as negative (outlier). The threshold value  $\tau$  changes with respect to the experiment (see Section 4).

## 4 Experiments

We devote this section to demonstrate the systematic workflow of our framework. Firstly, we pool images for a set of 21 object classes taken from the standard Caltech101<sup>2</sup>, using Flickr online photo sharing<sup>3</sup>. Each class contains 400 images

<sup>2</sup> [http://www.vision.caltech.edu/Image\\_Datasets/Caltech101](http://www.vision.caltech.edu/Image_Datasets/Caltech101)

<sup>3</sup> <https://www.flickr.com/>

as well as their corresponding tag lists ( $tag_i$ ). For simplicity, the number of crawled images has been defined in order to minimize the computational time of downloading images and managing their tags during the experiments. The effective number of classes have been normalized to 16, avoiding the classes that are composed by a bi-gram (i.e., two words). The remaining classes are: *accordion, bonsai, euphonium, face, laptop, menorah, nautilus, pagoda, panda, piano, pyramid, revolver, starfish, sunflower, umbrella, watch*. Since there are 400 images and 400 tag lists per class, the dataset is composed of 6400 images and 6400 tag lists.

To generate the ground-truth for our experiment in a more effective and efficient way, we build a graphical user interface (GUI) that allows us to manually label an image as positive or as negative with respect to the object class. For reliable manual classification, some guidelines are defined and adopted. If the following guidelines are satisfied, then an image is labeled as negative; otherwise as positive:

- an image is completely unrelated with the object specified by the category it belongs to;
- an image contains irrelevant parts of the object, that is, parts that alone are not sufficient to make the category object identifiable;
- an image contains only internal parts of the category object (like a cockpit of an airplane or an engine of a car);
- an image is a drawing or a caricature of the category object.

For each single feature we run several different experiments based on different strategies. Each strategy differs from the others with regard to the method used to compute the threshold. This produces different results in determining if a given tag list is associated to a positive or negative image.

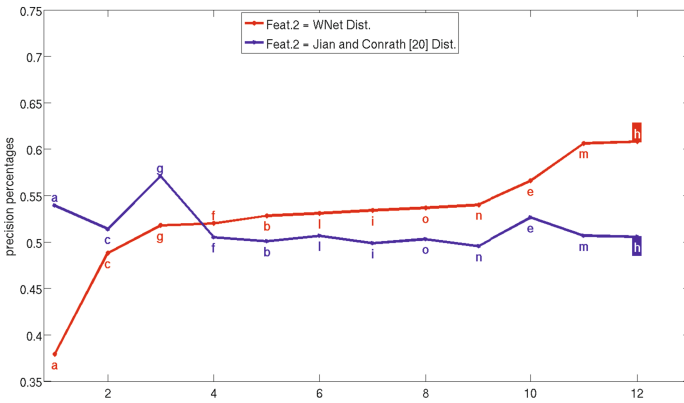
Referring to the algorithms described in the subsection 3.1, 3.2, 3.3, we give a brief explanation of the strategies associated to the threshold which produces the best discrimination results:

**Feature 1:** Based on experiments performances, we obtain the best result when searching if a keyword is found in the first three positions in the tag list. Surprisingly, this feature does not involve any cleaning mechanism of textual tags in the tag list (it avoids the step number 4 of algorithm 1). However, the feature takes the textual tags as they are provided by Flickr. At this point one may ask why using contaminated textual tags in a tag list is, unexpectedly, producing better results than the cleaned version. The answer lays in the “filtering” mechanism of the textual tags. Cleaning the tag list  $tag_i$  implies producing more single words ( $tag_j$ ) since the tag sentences are split. This increases the probabilities of finding the right match with the keyword, therefore a higher number of tags labeled as positive. This has been confirmed by the number of false positives generated using the other strategies, which is widely higher than the number of false positives produced by the strategy just described. To provide a better understanding of what happens if we do not perform any tag cleaning on the tag list, we present the following example: given the tag list relative to a negative



image of the *panda* class: “*zoo\_atlanta*”, “*taishan*”, “*giant\_panda*”, the keyword would not be matched since the substring matching is not performed. Therefore, the image is labeled as negative. This results change if we clean the tag list by splitting the sentences into single words. The cleaned tag list becomes: *zoo*, *atlanta*, *taishan*, *giant*, *panda*. In this case, the keyword would match with the 5<sup>th</sup> tag and therefore the image is now labeled as positive.

**Feature 2:** This feature uses two different measures: the standard semantic distance provided by WordNet, and the distance proposed by Jiang and Conrath in [20]. To select the one which produces the best results, we use both metrics to run the 12 strategies. We used these two distances since they are widely adopted in literature. The comparison results are shown in figure 2 .



**Fig. 2.** Summary results obtained by using WordNet and Jiang and Conrath distances in [20] in all the strategies. WordNet distance is outperforming in average in all of the strategies. We compute the precision rate for each strategy (*a, b, . . . , o*) as:  $\#TruePositive / (\#TruePositive + \#FalsePositive)$ .

Using WordNet distance as shown in figure 2, we observe constant increase in the average performances of all strategies. Therefore, in the following description we are mainly referring to the WordNet distance. The strategy based on the WordNet distance, which gives the best results, uses the following criteria to split the images set: defining the  $scores_i$  as the mean of the distances between the considered tags and the keyword:

$$scores_i = mean(Distance(tag_{(i,j)} - keyword))$$

$$Feat2(scores_i) = \begin{cases} positive & \text{if } mean(scores_i) \geq \tau \\ negative & \text{otherwise} \end{cases}$$

The best result is obtained using this strategy when the threshold is set to  $\tau = median(scores_I)$ , where the  $scores_I$  is the vector of all the  $scores_i$ .

**Feature 3:** The strategy based on the tag frequency feature, which produces the best results, compared with the other strategies, uses the following criteria to split the images set: defines the  $scores_i$  as the sum of the frequency values of the considered tags with respect to the keyword:

$$scores_i = \sum_{i=1}^{N_{images}} Fq(tag_{(i,j)})$$

$$Feat3(scores_i) = \begin{cases} positive & \text{if } mean(scores_i) \geq \tau \\ negative & \text{Otherwise} \end{cases}$$

We reach the best results when the threshold is set to  $\tau = mean(scores_I)$ , where the  $scores_I$  is the vector of all the  $scores_i$ .

## 5 Performances Evaluations

To assess the reliability of the experimental performances of the features described beforehand, we select  $n$  images labeled as positives from all the strategies and from Flickr. Hence, we count the true positives and the false positives that have been generated by the strategies and by Flickr (in this case, the false positives are the ones we manually label as negatives). Since the main goal of this framework is to generate a reliable dataset of images, for this reason, all of our strategies tend to produce more negative than positive labels. This behavior allows to minimize the number of the false positive labels generated during the experiments. Since not all strategies produce the same number of positive labels, to avoid the problem of getting some Null values, we fix  $n = \min(P - labels)$  of each feature. The selection of the  $n$  labels has been done randomly for Flickr, while for our strategies the first  $n$  are considered. To ensure the consistency of Flickr performances, we average the results produced after 10 random selections.

Table 1 displays the percentage values of the performances obtained using Flickr and our best strategies. The column  $\#P - labels$  contains the different  $n$  values used for each class. The column  $GT - Positives$  presents the number of true positives within the ground-truth.

To make the performances reported in the table more comparable, we recalculate the precision percentages by fixing  $n = 50$  positive labels<sup>4</sup> per class. Also in this case, the selection of the 50 labels has been done randomly for Flickr, while for our strategies it is referred to the first  $n$ . In figure 3, we provide the average values of each strategy for all the classes with  $n = 50$ .

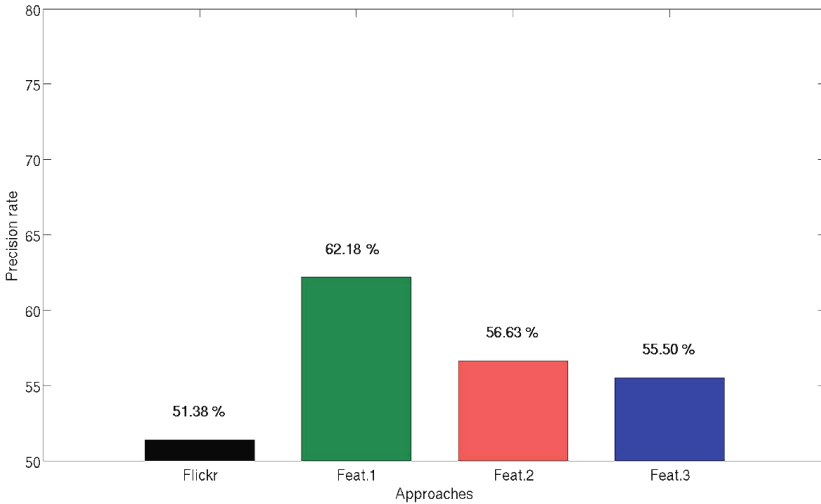
In this last case, an exception is done for the “*euphonium*” category, since it is composed by just 9 positive images also in the ground-truth.

At this point, one may be skeptical about the reliability of our strategies, since we are estimating their performances by considering only 50 images against the 400 downloaded. Therefore, if we observe how the performances change when we consider all the available positive labels shown in table 1, we are more confident

<sup>4</sup> This parameter has been set by considering the lowest common number of labels.

**Table 1.** Precision results obtained using all the features for 16 classes. Flickr provides the number of correct positive labels from the  $n$  images downloaded from the Flickr repository. *Feat* is an abbreviation for feature, where *Feat.1* refers to keyword position, *Feat.2* refers to semantic analysis, and *Feat.3* refers to tag frequency.

| Classes   | # P- labels | GT-Positives | Flickr | Feat.1 | Feat.2 | Feat.3 |
|-----------|-------------|--------------|--------|--------|--------|--------|
| watch     | 218         | 386 / 400    | 94.95  | 95.87  | 96.79  | 96.79  |
| sunflower | 178         | 379 / 400    | 93.26  | 97.19  | 96.63  | 96.63  |
| bonsai    | 119         | 362 / 400    | 90.76  | 90.76  | 92.44  | 88.24  |
| panda     | 182         | 359 / 400    | 89.56  | 90.11  | 32.31  | 97.25  |
| laptop    | 171         | 359 / 400    | 88.30  | 92.98  | 93.57  | 87.72  |
| pyramid   | 203         | 250 / 400    | 65.02  | 60.10  | 64.04  | 64.04  |
| starfish  | 170         | 211 / 400    | 49.41  | 60.00  | 56.47  | 53.53  |
| piano     | 50          | 105 / 400    | 37.50  | 58.33  | 37.50  | 70.83  |
| umbrella  | 175         | 164 / 400    | 37.14  | 41.14  | 41.71  | 44.00  |
| menorah   | 148         | 146 / 400    | 34.46  | 33.78  | 29.73  | 35.81  |
| accordion | 158         | 118 / 400    | 31.01  | 29.75  | 31.65  | 28.48  |
| pagoda    | 167         | 114 / 400    | 29.94  | 32.34  | 34.13  | 38.32  |
| face      | 135         | 120 / 400    | 28.15  | 31.11  | 25.19  | 27.41  |
| revolver  | 127         | 110 / 400    | 26.77  | 38.58  | 42.52  | 31.50  |
| nautilus  | 163         | 67 / 400     | 17.79  | 22.09  | 25.15  | 17.79  |
| euphonium | 8           | 9 / 400      | 0      | 62.5   | 0      | 0      |



**Fig. 3.** Summary of results of all the features by fixing  $n = 50$ . The highest precision is given using *feat.1* (i.e., keyword position).

on our results. Indeed, if we calculate the average of the positive labels considered in the last case, we can observe (see table 2) that the performances remain

constant when setting  $n \neq 50$ . The overall performance of our strategies still outperforms Flickr. In particular, using keyword position, the average performance obtained is encouragingly good (about 11% higher than Flickr). This information is further enriched since it provides us with a more reliable percentage value than the ones provided by the results of  $n = 50$ .

**Table 2.** The average performance of all the features when  $n = 50$  and  $n \neq 50$

| # P- labels | Flickr | Feat.1 | Feat.2 | Feat.3 |
|-------------|--------|--------|--------|--------|
| $\neq 50$   | 50.87  | 61.18  | 56.62  | 55.50  |
| $= 50$      | 50.87  | 62.18  | 56.62  | 55.50  |

## 6 Conclusions

We have presented a framework to support the entire process of image and textual retrieval from search engines that, given an input keyword, performs a statistical and a semantic analysis and automatically builds a training set. We have conducted several experiments to validate our assumptions about the analysis of textual information and the evaluation that we have provided on three investigated methods have shown that the position of tags, their order, is relevant. We have investigated the semantic aspects by using semantic distance. Unfortunately, the results achieved show modest benefit for the adopted semantic features. However, the methods suggested are currently under continuous experimentation and need to be further investigated. In particular, we consider for future work to explore the use of different search engines such as Google<sup>5</sup>, ImageNet<sup>6</sup>, InstaGram<sup>7</sup> or Pinterest<sup>8</sup> to check if they are interchangeable or can be combined to improve performances. We plan also to extend and investigate other semantic features related to ontological relationships of textual information and combine them with the aim of creating a waterfall model which combines different strategies.

**Acknowledgements.** This research was supported by the VISCOSO project financed by the Autonomous Province of Trento through the “Team 2011” funding programme.

## References

1. Helmer, S., Meger, D., Viswanathan, P., McCann, S., Dockrey, M., Fazli, P., Southey, T., Muja, M., Joya, M., Jim, L., Lowe, D.G., Mackworth, A.K.: Semantic

<sup>5</sup> <http://www.google.com>

<sup>6</sup> <http://www.image-net.org>

<sup>7</sup> <http://instagram.com/>

<sup>8</sup> <http://www.pinterest.com>

- robot vision challenge: current state and future directions. In: IJCAI workshop (2009)
2. Cheng, D.S., Setti, F., Zeni, N., Ferrario, R., Cristani, M.: Semantically-driven automatic creation of training sets for object recognition. *Computer Vision and Image Understanding* 131, 56–71 (2014)
  3. Lew, M.S., Sebe, N., Djeraba, C., Jain, R.: Content-based multimedia information retrieval: State of the art and challenges. *ACM Trans. Multimedia Comput. Commun. Appl.* **2**(1), 1–19 (2006)
  4. Liu, Y., Xu, D., Tsang, I.W., Luo, J.: Textual query of personal photos facilitated by large-scale web data. *IEEE Trans. Pattern Anal. Mach. Intell.* **33**(5), 1022–1036 (2011)
  5. Heymann, P., Paepcke, A., Garcia-Molina, H.: Tagging human knowledge. In: Proceedings of the Third ACM International Conference on Web Search and Data Mining. WSDM 2010, pp. 51–60. ACM, New York (2010)
  6. Lin, Z., Ding, G., Hu, M., Wang, J., Ye, X.: Image tag completion via image-specific and tag-specific linear sparse reconstructions. In: 2013 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1618–1625, June 2013
  7. Wu, L., Jin, R., Jain, A.K.: Tag completion for image retrieval. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(3), 716–727 (2013)
  8. Fei-Fei, L., Perona, P.: A bayesian hierarchical model for learning natural scene categories. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2005. CVPR 2005, vol. 2, pp. 524–531 (2005)
  9. Fei-Fei, L., Fergus, R., Perona, P.: Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. *Computer Vision and Image Understanding* **106**(1), 59–70 (2007)
  10. Li, L.J., Li, F.F.: Optimol: Automatic online picture collection via incremental model learning. *International Journal of Computer Vision* **88**(2), 147–168 (2010)
  11. Fergus, R., Fei-Fei, L., Perona, P., Zisserman, A.: Learning object categories from google’s image search. In: Tenth IEEE International Conference on Computer Vision, 2005. ICCV 2005. vol. 2, pp. 1816–1823 (2005)
  12. Gilbert, A., Bowden, R.: A picture is worth a thousand tags: automatic web based image tag expansion. In: Lee, K.M., Matsushita, Y., Rehg, J.M., Hu, Z. (eds.) ACCV 2012, Part II. LNCS, vol. 7725, pp. 447–460. Springer, Heidelberg (2013)
  13. Dubinko, M., Kumar, R., Magnani, J., Novak, J., Raghavan, P., Tomkins, A.: Visualizing tags over time. In: Proceedings of the 15th International Conference on World Wide Web. WWW 2006, pp. 193–202. ACM, New York (2006)
  14. Ahern, S., Naaman, M., Nair, R., Yang, J.: World explorer: Visualizing aggregate data from unstructured text in geo-referenced collections. In: Proceedings of the Seventh ACM/IEEE-CS Joint Conference on Digital Libraries, pp. 1–10. ACM Press (2007)
  15. Spain, M., Perona, P.: Some objects are more equal than others: measuring and predicting importance. In: Forsyth, D., Torr, P., Zisserman, A. (eds.) ECCV 2008, Part I. LNCS, vol. 5302, pp. 523–536. Springer, Heidelberg (2008)
  16. Ames, M., Naaman, M.: Why we tag: motivations for annotation in mobile and online media. In: CHI 2007: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pp. 971–980. ACM Press, New York (2007)
  17. Wu, L., Hua, X.S., Yu, N., Ma, W.Y., Li, S.: Flickr distance: A relationship measure for visual concepts. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **34**(5), 863–875 (2012)

18. Hwang, S.J., Grauman, K.: Learning the relative importance of objects from tagged images for retrieval and cross-modal search. *Int. J. Comput. Vision* **100**(2), 134–153 (2012)
19. Vijayanarasimhan, S., Grauman, K.: Keywords to visual categories: Multiple-instance learning for weakly supervised object categorization. In: *IEEE Conference on Computer Vision and Pattern Recognition, 2008. CVPR 2008*, pp. 1–8 (June 2008)
20. Jiang, J.J., Conrath, D.W.: Semantic similarity based on corpus statistics and lexical taxonomy. *CoRR cmp-lg/9709008* (1997)
21. Fellbaum, C.: *WordNet: An Electronic Lexical Database*. Language, Speech and Communication. Mit Press (1998)