

Pyramid Histograms of Visual Words (PHOW)

Natalia Andrea Durán Castro
Universidad de los Andes
na.duran@uniandes.edu.co

Ana M. Velosa Orduz
Universidad de los Andes
am.velosa@uniandes.edu.co

Abstract

1. Introduction

1.1. Datasets

Imagenet is a dataset of images that are organized according to a hierarchy presented on WordNet. This is a large database of lexical in English. It is grouped depending of conceptual-semantic and lexical relations [4]. The Imagenet dataset categories consists of nouns in the synonym set as per WordNet. There are more than 80,000 synonym sets of nouns in WordNet [4]. For each category, the ImageNet dataset provides 10,000 images with human annotations and quality control [3].

Caltech 101 is a dataset of 101 different categories. In each category, there are between 40 to 800 images. The size of each image is around 300x200 pixels. Also, this dataset included anotations. According to the error reports of Caltech, the performance is better when all the images of test are used and, it is not normalized to all the categories [2].

1.2. PHOW

Pyramid Histogram of Visual Words (PHOW) is an unsupervised trained model based on the engineered models to use in computer vision [1]. This is an extension of the bag of words model (BoW), where the images are considered as words. In this way, a BoW model considers the images features as words. Then, the bag words would be a sparse vector(vector where the majority of values are zero) with these words [1]. Thereby, it is possible to know the frequency in which a word(image feature) is repeated on the image.

However, the BoW method does not consider the spatial information. For this reason, the PHOW method subdivided the image into incrasiling sub-regions as pyramids. The histogram is computed in each one of these regions.

The scale invariant feature transform (SIFT) is an algorithm to detect and describe local features in images [1].

data set from Caltech.jpg

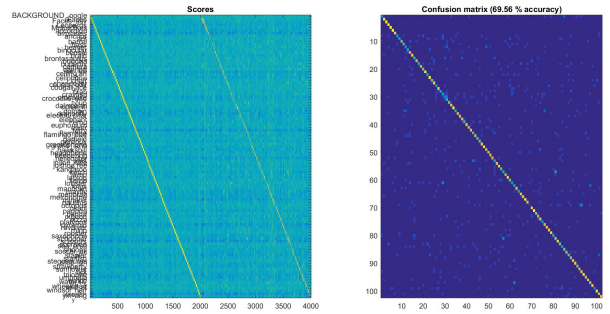


Figure 1. ACA of the complete Caltech Dataset without modifications

The difference between this algorithm and PHOW is that SIFT does not read the image features as words. Then, the PHOW can use the SIFT descriptors with k-means clustering to transform them in a visual vocabulary that can be used as words. Due to SIFT is invariant to changes in scale and orientation of the images, the PHOW would be invariant to this.

The difference between textons and PHOW is that textons just considers the spatial information on a color space of representation. On the other side, PHOW takes account more information of the image features. This shows better results in large-scale datasets.

1.3. Methodology

The hyperparameters used in PHOW are

2. Results and discussion

2.1. Caltech

In the figure 1 we can see the value of the ACA that was obtained with the Caltech database with all the images. On the other hand, in the blah image we can observe the value that the ACA takes with the imageNet200 database as if it were a small problem, in which a smaller amount of class is taken.

We can see that the accuracy of imageNet200 is much

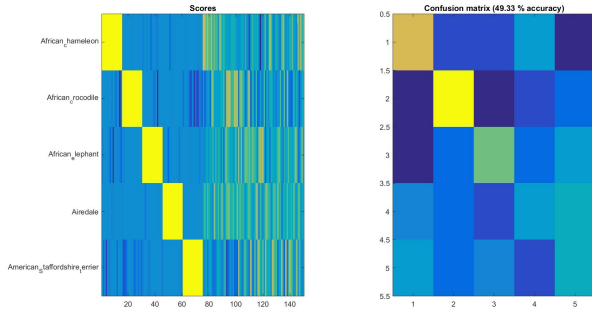


Figure 2. ACA of the tiny problem of the imageNet200 DataSet

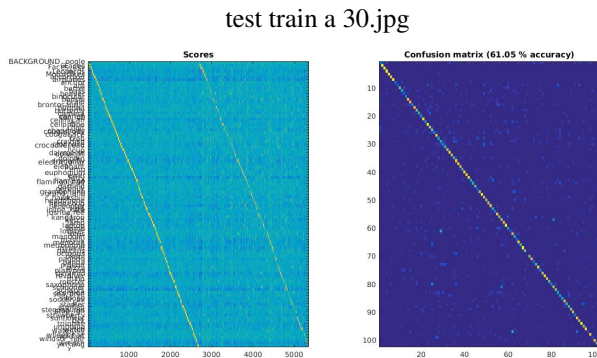


Figure 3. ACA of the complete Caltech Dataset with resize of the amount of images in train and test to 30

lower than that of Caltech which is expected because having a greater number of classes the algorithm should acquire a higher level of sensitivity and therefore be able to classify better. The parameters we changed

The parameters that were changed each time the algorithm was run were the number of images that were used for validation and for training, the value of the step and finally, what value was given to the C parameter of the SVM. By modifying the value only of the step we could see that the value of the ACA did not change no matter how big or small we put the error range of C. On the other hand, when changes were made with respect to the step, we observed the same as with parameter C, which did not significantly change the ACA

Finally, when the number of images in the validation and the interweaving was modified, it could be observed that for more images and for less images the value of the ACA diminished, as can be seen in 3 and 4

Taking into account all of the above, it is important to emphasize that a change of 20 images in the validation and test images generates a very significant change since it goes from an ACA 61.05 to an ACA of 1.37 percent.

The num patialX / Y was also modified to dobre, that is, [4 8] but no change was observed in the aca obtained as observed in ??

size of train and test to 10.jpg

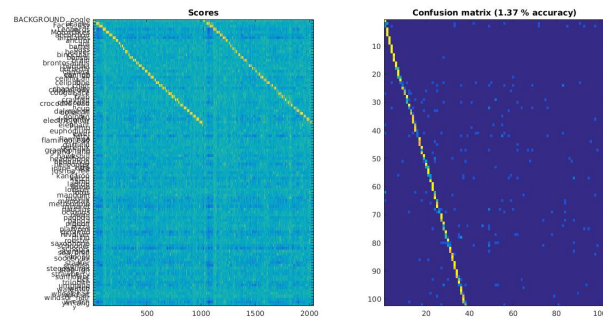


Figure 4. ACA of the complete Caltech Dataset with resize of the amount of images in train and test to 10

on the numSpatialX/Y to [4 8].jpg

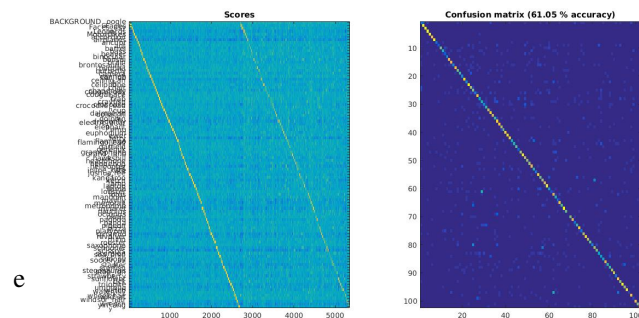


Figure 5. ACA of the complete Caltech Dataset with resize of the numSpatialX/Y to [4 8]

References

- [1] S.-M. Khaligh-Razavi. Mrc cognition and brain sciences unit, cambridge university, cambridge, uk. pages --. 1
- [2] R. F. L. Fei-Fei and P. Perona. One-shot learning of object categories. IEEE Trans. Pattern Recognition and Machine Intelligence. In press. Retrieved from: http://www.vision.caltech.edu/Image_Datasets/Caltech101/Description.1
- [3] O. R. J. K. J. D. A. B. Li Fei-Fei, Kai Li. Imagenet. Retrieved from: image-net.org/about-people.1
- [4] G. A. Miller. Wordnet. Retrieved from: word-net.princeton.edu/people.1