

# Large-scale Image Labeling via Mapreduce Topic Modeling

[Spring 2012 CMSC828G Final Project Report] \*

Khoa Doan<sup>†</sup>  
University of Maryland  
Department of Computer  
Science  
College Park, MD  
trovato@corporation.com

Rahmatri Mardiko<sup>‡</sup>  
University of Maryland  
College Park, MD 20742  
mardiko@cs.umd.edu

Ang Li<sup>§</sup>  
University of Maryland  
Depart. of Computer Science  
College Park, MD  
angli@cs.umd.edu

## ABSTRACT

Large scale computer vision is recently being popular due to the emerging of large scale imagery data and in need of more generally trained visual models. Recent advances in feature extraction provide a feasible and succinct way to represent image regions. However, most of these features are computationally heavy to extract. Since there is no explicit method to essentially speed up feature extraction, parallelization becomes the most comfortable strategy for this task. Due to the inherent high dimensionality of visual data, extracted features can be very noisy and thus not representative for images. Feature quantization is introduced to group similar features into the same low level semantics. K-means clustering is one of the usual choices for quantization while it may suffer from either time or space problems in ordinary environment. Latent Dirichlet Allocation is one of the techniques to discover higher level topic semantics for the context, which was popular in natural language processing. In this project, we unearth the potential of adapting computer vision tasks such as low level feature extraction and higher level image understanding into the MapReduce framework for large scale image dataset. We adopt the MapReduce LDA (Mr.LDA) to find topics across the images and discuss its possible extensions with respect to the image domain.

## Keywords

Topic modeling, image labeling, parallelization, MapReduce framework

## 1. INTRODUCTION

\*Data-Intensive Computing with MapReduce instructed by Prof. Jimmy Lin

<sup>†</sup>Dr. Trovato insisted his name be first.

<sup>‡</sup>The secretary disavows any knowledge of this author's actions.

<sup>§</sup>This author is the one who did all the really hard work.

Large scale visual data analysis and machine learning have been recently received much attention during the past few years. Traditional computer vision research focuses on small datasets of images or videos which makes the generalization of these methods difficult. In the current world of big data, a lot of imagery data have been present in the Internet. How to make use of the large scale of data to explore a better visual model has been one of the central problems in the current community. However, one natural problem arised from this task is heavy computational load. Computer vision tasks are usually involved with feature engineering which requires a lot of time and space for experiments. Fortunately, the MapReduce framework provides a reliable approach to large scale data processing.

In this work, we explore the potentials of using MapReduce framework for a standard computer vision task i.e. image labeling. The objective of the image labeling task is to find underlying semantics for each of the image pixels and to find groups of pixels that belong to the same object or semantic. Image labeling has been investigated for decades, although the scale of the dataset is limited due to computational issues. However, the demand of large scale image labeling turns out to be more and more clear in the recent years. One of the reasons is that people are being aware of using computer techniques to assist human annotation for specific imagery data such as remote sensing data. Millions of satellite images are generated for every day and the task of understanding these data is never feasible for human to do exhaustively. The automatic way will benefit the community tremendously in different areas such as surveillance, city planning, national defense, etc.

The rest of this paper is organized as follows. Section 2 briefly introduces the system structure of this project. MapReduce framework for visual feature extraction is discussed in section 3. Details of Mr.LDA is introduced in section 4. In section 6, a few implementation details are discussed. Section 7 shows the evaluation methods and experimental results for this project. Discussion on possible extensions from LDA to Spatial LDA is presented in section 8. Finally, the paper concludes in section 9.

## 2. SYSTEM OVERVIEW

## 3. FEATURE EXTRACTION

Feature extraction often dominates the most computational resources in computer vision tasks. In this section, we introduce a MapReduce based framework to extract visual features efficiently which distributes the computation loads for raw feature extraction and scales up the construction of feature codebook via MapReduce K-means clustering.

### 3.1 Scale Invariant Feature Transform (SIFT)

In this work, the Scale Invariant Feature Transform (SIFT) [4] is adopted to describe local regions across the images. Generally, a SIFT descriptor represent a region bounded by a rectangular box centered at  $(x_c, y_c)$ . The bounding box is uniformly divided into  $K \times K$  parts. The gradients are computed in these sub-regions and the gradient magnitudes in  $N_b$  (the number of bins) discretized orientations are computed. A histogram is then constructed to represent the distribution of the gradient magnitudes. Therefore, for each region, the SIFT descriptor is a  $K^2 N_b$  dimensional histogram vector.

One of the advantages of SIFT descriptors is its invariance to image rotation and translation, which makes this feature very popular among the computer vision community. We base our project on the MPI-CBG JavaSIFT library[2].

### 3.2 Dense Sampling SIFT Features

In order to apply SIFT features to represent the whole image for the task of image labeling, each of the pixels ideally should be described. Due to the fact that a neighborhood of pixels usually have little difference and belong to the same semantics, we sample series of small rectangular regions with a parameter STEPSIZE densely across each of the image and compute the SIFT description for each of the regions. The STEPSIZE controls the number of pixels between two consecutive centers of bounding boxes. Thereafter, each image can be converted from visual data to a list of (KEY,VALUE) pairs where the KEY is the index  $(g_i, x_i, y_i)$  of the regions and the VALUE is the SIFT feature vector  $f_i$ .  $g_i$  is the id of the image that contains the  $i$ -th region and  $(x_i, y_i)$  is the center location of the  $i$ -th region.

### 3.3 Image Input Aggregation

Given a large number of image files as input, there are several ways of loading them into a MapReduce job. One possible way is creating a text file containing a list of file names and let the mapper reads the files from HDFS and processes them while the reducer does nothing. An alternative way is similar to the first except the reading and processing is moved to the reducer while the mapper just passes the (key,value) pairs. Both are relatively simple and easy to implement. However, due to the small-files problem in Hadoop[5], neither the first and the second approach achieves the best performance in terms of processing speed.

Since Hadoop framework works best with large files, we need to aggregate the images in the dataset into a few big files. This step allows us to speed-up the feature extraction (and other images processing tasks) given a large number of image files as input. We adopt the approach presented in[6] which creates a MapReduce job that packs multiple files into a `SequenceFile`. This aggregation is performed as a preprocessing step for the other jobs that take images as input.

## 3.4 MapReduce-based Feature Extraction

A list of image paths is generated in order to locate the image files. Each of the paths is assigned a number as image ID. The MapReduce job takes the path list as input and output sequence files containing region indices as keys and SIFT feature descriptors as values.

---

**Algorithm 1:** MAX finds the maximum number

---

**Input:** A finite set  $A = \{a_1, a_2, \dots, a_n\}$  of integers

**Output:** The largest element in the set

```

 $max \leftarrow a_1;$ 
for  $i \leftarrow 2$  to  $n$  do
    if  $a_i > max$  then
         $max \leftarrow a_i;$ 
return  $max;$ 

```

---

### 3.5 Building Feature Codebook

The raw features extracted from above generally have two problems in representing the images. On the one hand, the dimension of each feature vector is typically 128. For the similarity measure of each pair of small regions, at least 128 times of multiplications are necessary for Euclidean distances. This makes the further processing intractable even using parallelization framework. On the other hand, the feature vector is a histogram of oriented gradients which contains a lot of noises. Therefore, feature quantization is introduced to further processing the features and construct a "codebook" for the features. In the codebook, nearby feature vectors are grouped into the same index because they should belong to the same visual semantics.

K-means clustering is usually the choice for doing feature quantization. However, the clustering takes more features and thus needs much more memory spaces in order to exhaustively compute the distances between cluster centers to each of the points, as the scale of data becomes larger. The MapReduce framework generally resolve this time and space problem because not only the computation is distributed into different nodes but also the intermediate results of distances are stored almost uniformly among all of the nodes. In our project, we adopt Mahout K-means clustering [1] for building the codebook in Hadoop environment.

## 4. MAPREDUCE LDA

[7]

## 5. MAPREDUCE SPATIAL LDA

## 6. IMPLEMENTATION REMARKS

In this section we present some implementation details that we did in the project.

### 6.1 Images to Pixels Conversion

In the evaluation phase of the project we need to compare the LDA output and the true labels. To enable pixel-by-pixel evaluation we extract each pixel in the ground truth images and store them as (key,value) pairs. The key, which contains the image id and the pixel position  $(x,y)$ , is of type `TripleOfInts` whereas the value is of type `IntWritable` and it contains the pixel value. This images-to-pixels conversion can actually be performed as map-only MapReduce job.

However, sometimes it is useful to group pixels that have the same semantic together. Here we can apply value-to-key conversion so the pixels that have the same labels are grouped together when they arrive at the reducer.

---

**Algorithm 2:** MapReduce Convert Images to Pixels

---

```

mapperIntWritable key, BytesWritable value image =
read(value); width,height = size(image); for
y = 1 → height do
-
x = 1 → width pixel = getpixel(image,x,y); emit
(pixel, (key,x,y));

reducerIntWritable key, Iterable<TripleOfInts> values
iter = iterator(values) while iter.hasNext() do
-
emit (iter.next(), key);

```

---

## 6.2 Pixels to Images Conversion

Also for the evaluation purpose, we need to build images from a set of pixels. This task is the inverse of the previous task. The mapper takes ((image\_id,x,y),pixel) pairs and produce (image\_id,(x,y,pixel)) as intermediate key value pairs so the pixels that belong to the same image are grouped together in the reducer. In the reducer the pixels are used to build the image object and convert it to BytesWritable. Since the output is in the form of SequenceFile, we still need to read the output separately to get the individual image files.

---

**Algorithm 3:** MapReduce Convert Pixels to Images

---

```

mapperTripleOfInts key, IntWritable value image_id,x,y
= getmembers(key); emit (image_id, (x,y,value));

reducerIntWritable key, Iterable<TripleOfInts> values
image = new image(width,height); iter =
iterator(values) while iter.hasNext() do
-
x,y,pixel = getmembers(value);
setpixel(image,x,y,pixel);
imseq = getbytes(image); emit (key, imseq);

```

---

## 7. EXPERIMENTAL EVALUATION

### 7.1 Dataset

#### 7.1.1 MSRC Image Labeling Dataset

The MSRC image labeling dataset (version 1) [3] contains 240 images and 9 object classes in total. Each of the images is pixel-wise labeled. Fig.(1) shows some sample images from MSRC dataset.



**Figure 1:** Sample images from MSRC Image Labeling Dataset

#### 7.1.2 Eastern Coast Satellite Image Dataset

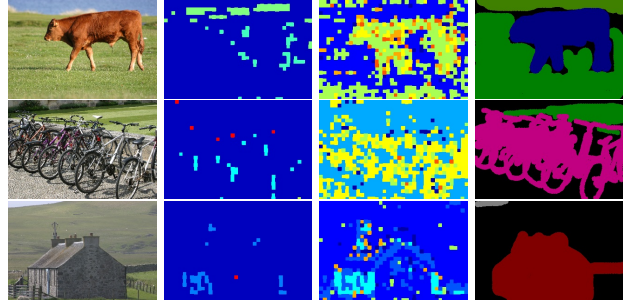
We also collect satellite images using Google Maps APIs along the eastern coast in the United States. Each of the image is of 8-bit colormap PNG format and contains  $402 \times 415$  pixels. Due to the capacity of our cluster, we pick 10,000 images amongst the dataset for the experiments. Fig.(2) shows some sample images from Satellite image dataset.



**Figure 2:** Sample images from Eastern Coast Satellite Image Dataset

### 7.2 Qualitative Evaluation

Images with labels in different colors are recovered from the output of LDA to show qualitatively which parts belong to the same semantics. Figure 3 presents three sample images and the output labels with 9 and 14 topics. The ground truth consists of 14 distinct semantics. However, the performance of LDA is significantly higher with 9 topics than 14 topics.



**Figure 3:** Sample LDA output from MSRC dataset. The second and the third column show the output of labeling with LDA using 14 and 9 topics. The last column show the ground truth images.

### 7.3 Quantitative Evaluation

Traditional computer vision dataset has ground-truth for evaluation. For the MSRC image labeling data, each image is corresponded to a groundtruth image in which the whole image is divided into several components. Each component is annotated using a unique RGB color. In our output, each pixel is assigned to a type of semantics. Since our approach is completely unsupervised, the truth meaning of each semantic is not understandable. One possible way to evaluate the performance is to enumerate every possible correspondence between the two set of semantics. Apparently this strategy only works when the number of topics is small enough. It becomes difficult to construct a one-one correspondence between our labels and the groundtruth labels. In fact, how to accurately evaluate the performance of image segmentation tasks is still an open problem.

In this project, we propose to evaluate the results for each of the topics separately. For each semantic in the ground truth, a binary image can be constructed, regions with value 1 are those containing such semantic and otherwise not. For the output of the resulted topics, the topic that has the most number of occurrences in the semantic regions are chosen as a correspondence. The precision and recall are then computed according to that topic.

$$\text{precision} = \frac{\# \text{TruePositive}}{\# \text{TruePositive} + \# \text{FalsePositive}} \quad (1)$$

$$\text{recall} = \frac{\# \text{TruePositive}}{\# \text{TruePositive} + \# \text{FalseNegative}} \quad (2)$$

## 8. DISCUSSION: SPATIAL LATENT DIRICHLET ALLOCATION

## 9. CONCLUSION

## 10. ACKNOWLEDGMENTS

The authors would like to thank Prof. Jordan Boyd-Graber and Prof. Jimmy Lin for their consistent help in this project.

## 11. REFERENCES

- [1] Mahout: <http://mahout.apache.org/>.
- [2] mpi-cbg javasift library:  
<http://fly.mpi-cbg.de/saalfeld/projects/javasift.html>.
- [3] Msrc dataset: <http://research.microsoft.com/en-us/projects/objectclassrecognition/>.
- [4] D. G. Lowe. Distinctive image features from Scale-Invariant keypoints. *Int. J. Comput. Vision*, 60(2):91–110, Nov. 2004.
- [5] T. White. The small files problem:  
<http://blog.cloudera.com/blog/2009/02/the-small-files-problem/>.
- [6] T. White. *Hadoop: The Definitive Guide*, pages 239–245. O’Reilly Media, Inc., 1st edition, 2009.
- [7] K. Zhai, J. B. Graber, N. Asadi, and M. L. Alkhoulja. Mr. LDA: a flexible large scale topic modeling package using variational inference in MapReduce. In *Proceedings of the 21st international conference on World Wide Web, WWW ’12*, pages 879–888, New York, NY, USA, 2012. ACM.

## APPENDIX