



UNIVERSITY OF INFORMATION  
TECHNOLOGY

COMPUTER VISION PROJECT

**Objects Classification on Caltech 101**

*Truong Phuc Anh*

*Lam Han Vuong*

Supervised by

Mr.Ngo Duc Thanh

2 January 2018

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Dataset</b>	<b>2</b>
<b>3</b>	<b>Methods</b>	<b>2</b>
3.1	Flow . . . . .	2
3.2	Bag Of Visual Words (BOV) . . . . .	3
3.3	K-Neareast Neighbor (KNN) . . . . .	6
3.4	Support Vector Machines (SVMs) . . . . .	7
<b>4</b>	<b>Experiments</b>	<b>12</b>
4.1	Data preprocessing . . . . .	12
4.2	Feature extraction . . . . .	12
4.3	KNN test . . . . .	12
4.4	SVMs test . . . . .	12
<b>5</b>	<b>Conclusions</b>	<b>14</b>
<b>6</b>	<b>References</b>	<b>15</b>
<b>7</b>	<b>Appendices</b>	<b>15</b>

## List of Figures

1	Object recognition challenges example . . . . .	1
2	Generic flow for object recognition . . . . .	2
3	BOV-SVM flow chart . . . . .	3
4	BOV example . . . . .	4
5	KNN example . . . . .	6
6	Compare hyperplanes in SVMs . . . . .	9
7	The margin example . . . . .	10

## List of Tables

1	SVM result with linear kernel. . . . .	13
2	SVM result with poly kernel. . . . .	13
3	SVM result with rbf kernel. . . . .	14

# 1 Introduction

Object recognition is a common task of computer vision for identifying a specific object in a digital image or video. Humans can recognize multiple objects in images with little effort, but this task is still a big challenge for computer vision systems due to the fact that the objects may vary somewhat in different view points, in many different sizes, scales, types or even when they are translated or rotated (in video).



Figure 1: Chairs with different view points, sizes, materials make the recognition task more challenging.

Object recognition is useful in applications such as video stabilization [1], Manufacturing Quality Control [2], Automated vehicle parking systems [3], etc.

Many approaches to the task have been implemented over the years. Common techniques include deep learning based approaches such as convolutional neural networks, and feature-based approaches using edges, gradients, histogram of oriented gradients (HOG), Haar wavelets, and linear binary patterns.

In this project, we end up with a simple feature-based approach using Bag of words (BOW) with two different learning algorithms K-Nearest Neighbor (KNN)

and Support Vector Machine (SVM) on Caltech101 just for practising what we have learn from the class.

## 2 Dataset

We use Caltech 101 for this project. Pictures of objects belonging to 101 categories. About 40 to 800 images per category. Most categories have about 50 images. Collected in September 2003 by Fei-Fei Li, Marco Andreetto, and Marc 'Aurelio Ranzato. The size of each image is roughly 300 x 200 pixels.

For more infomation about the dataset, please visit [http://www.vision.caltech.edu/Image\\_Datasets/Caltech101/Caltech101.html](http://www.vision.caltech.edu/Image_Datasets/Caltech101/Caltech101.html).

## 3 Methods

### 3.1 Flow

We used the generic flow for this recognition task with *Bag of visual words* as feature extractor and *KNN*, *SVM* as learning algorithms.

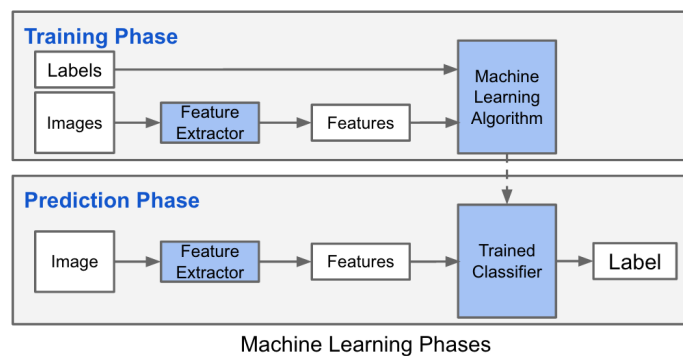


Figure 2: Generic flow for object recognition

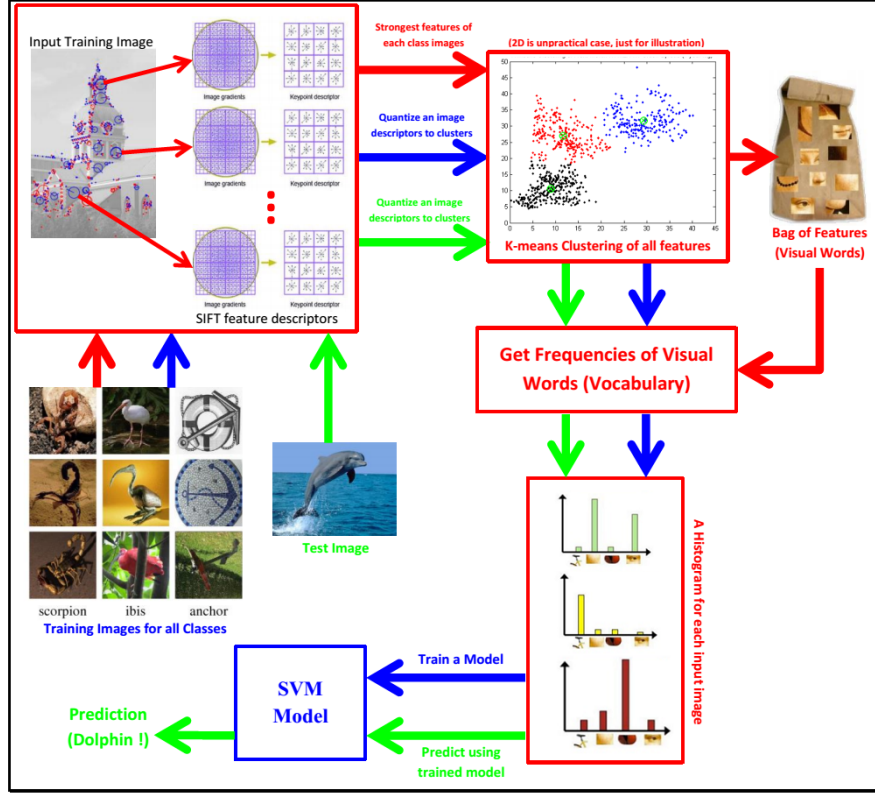


Figure 3: The detail process of recognition with BOW and SVM

## 3.2 Bag Of Visual Words (BOV)

### Over view

The bag-of-words model is a [vector space model](#) representing text data when modeling text with machine learning algorithms. In practice, this model is mainly used as a tool of feature generation.

### Basic idea

In document classification after transforming the text into a "bag of words", we can calculate various measures (also call features) to characterize the text. The most common type of characteristic is the histogram of words (based on the number of times a word appears in the document).

In order to use bag of words (from text) in computer vision (for images), we can treat image features as words. Thus, we have *Bag Of Visual Words*. The vocabulary can be built up by clustering all feature points collected from a suitable source (Ex. the training set) with each cluster represents for a word. After we have the vocabulary, the histogram of words will be used as feature vector of an image.

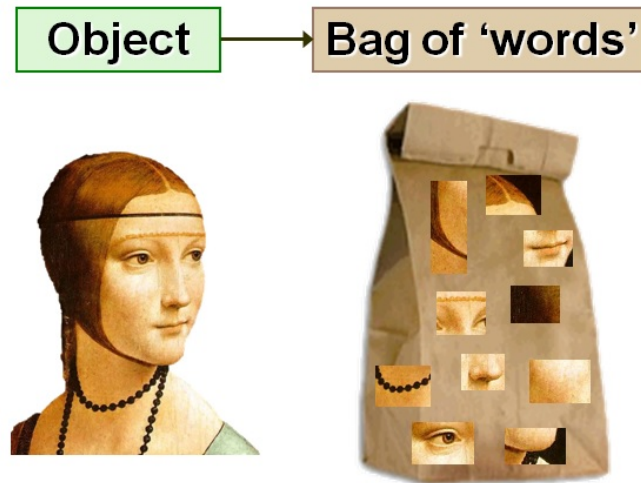


Figure 4: An example for BOV

In this project, we use SIFT as feature extractor and K-Means as clustering algorithm to build the vocabulary.

### Pseudo-code

```
1 # Create a vocabulary from a set of images
2 def create_vocabulary(images):
3     # 1. Get key points (interested points) from all images
4     key_points = get_key_points(images)
5
6     # 2. Calculate descriptors for those key points.
7     descriptors = get_descriptors(images, key_points)
8
```

```

9     # 3. Run K-Means clustering on those descriptor
10    cluster_centers = run_kmeans_on(descriptors)
11
12    # 4. Each cluster center is a word in vocabulary
13    return cluster_centers as vocabulary

```

---

```

1  # Calculate histogram of visual words (as a feature vector) for an image
2  def cal_histogram(image, vocabulary):
3      # 1. Create an empty histogram
4      n_words = len(vocabulary)
5      histogram = np.zeros(n_words) # zeros array with length = n_words
6
7      # 2. Calculate all descriptors in this image
8      key_points = get_key_points(image)
9      descriptors = get_descriptors(key_points)
10
11     # 3. For each descriptor, find out which word it is (in vocabulary).
12     for descriptor in descriptors:
13         word_index = which_word(descriptor, vocabulary)
14         histogram[word_index]++ # increase histogram bin
15
16     return histogram

```

---

```

1  # Find out which word an descriptor is (in vocabulary)
2  def which_word(descriptor, vocabulary):
3      n_words = len(vocabulary)
4      dist = np.zeros(n_words) # zeros array with length = n_words
5      for i in range(n_words)
6          dist[i] = L1_distance(vocabulary[i], descriptor)
7      return min_index_of(dist)

```

---



### 3.3 K-Neareast Neighbor (KNN)

#### Over view

KNN is a type of [instance-based learning](#), or [lazy learning](#) which means we don't need to train a real model for redicting.

#### Basic idea

The basic idea is with any unknown sample  $x$  (or a query point), finding  $k$  *nearest samples* in training set and assigning the label which is most frequent among  $k$ -nearest *samples* to that query point.

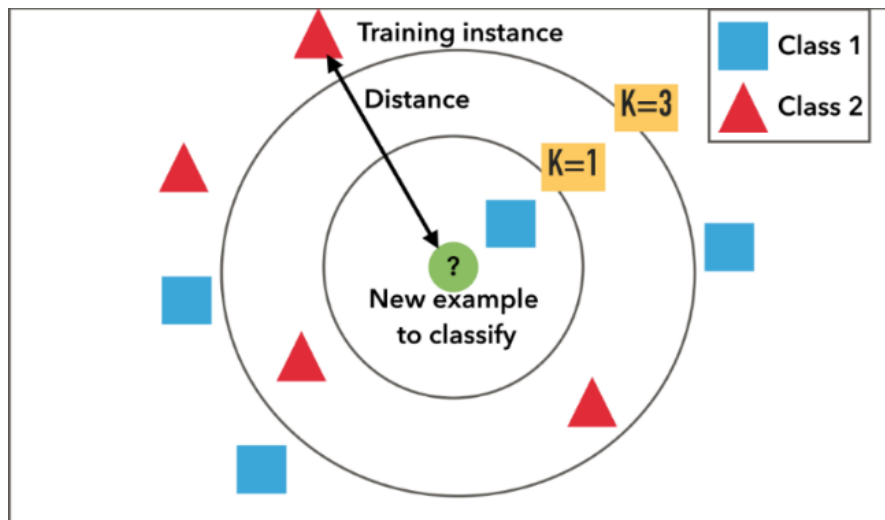


Figure 5: An example of KNN with different K

A common method used for calculate the distance metric between samples is [Euclidean distance](#). We also use this model in this project for a simple implementation.

#### Pseudo-code

---

```

1 # Find the label of an unknown sample x
2 # x: unknown sample
3 # X: training data
4 # Y: labels of training data
5 # k: parameter k
6 def knn(x, X, Y, k):
7     n_samples = length(X)
8
9     # 1. Calculate the distance between x and all training samples in X
10    for i in range(0, n_samples):
11        distances[i] = Euclidean_distance(x, X[i])
12
13    # 2. Get index of k-th minimum distances
14    k_indeces = get_minimum_indeces(k, distances)
15
16    # 3. Get voted labels from k_indeces
17    voted_labels = labels[k_indeces]
18
19    # 4. return the most frequent among voted_labels as the label of x
20    return most_frequent(voted_labels)

```

---

### 3.4 Support Vector Machines (SVMs)

#### Over view

SVMs are [supervised learning](#) models with associated learning algorithms that analyze data used for classification and regression analysis. [4]. In this project we only consider about classification.

In 1957, a simple linear model called the Perceptron was invented by Frank Rosenblatt to do classification.

A few years later, Vapnik and Chervonenkis, proposed another model called the "Maximal Margin Classifier", the SVM was born.

Then, in 1992, Vapnik et al. had the idea to apply what is called the Kernel Trick, which allow to use the SVM to classify linearly nonseparable data.

Eventually, in 1995, Cortes and Vapnik introduced the Soft Margin Classifier which allows us to accept some misclassifications when using a SVM.

So just when we talk about classification there is already four different Support Vector Machines:

1. The original one : the Maximal Margin Classifiers
2. The kernelized version using the Kernel Trick
3. The soft-margin version
4. The soft-margin kernelized version (which combine 1, 2 and 3)

For more detail about the history, you can read [here](#).

In this project, we implement the last version for multi-classification with two approaches *one vs one* and *one vs rest*

## Basic idea

In the case of SVM, it learns a linear model (a hyperplane) to separate data. Let's take a quick look at how the first SVM (also call hard margin) is built up in simple case (linear separable data).

Before SVM, Perceptron Learning Algorithm is a simple model for finding a hyperplane to separate data. However, its biggest weakness is that it will not find the same hyperplane every time but not the best one (because of the randomness, you can find the detail [here](#)). SVMs solve this problem by finding the *best hyperplane* (we call this the **optimal hyperplane**).

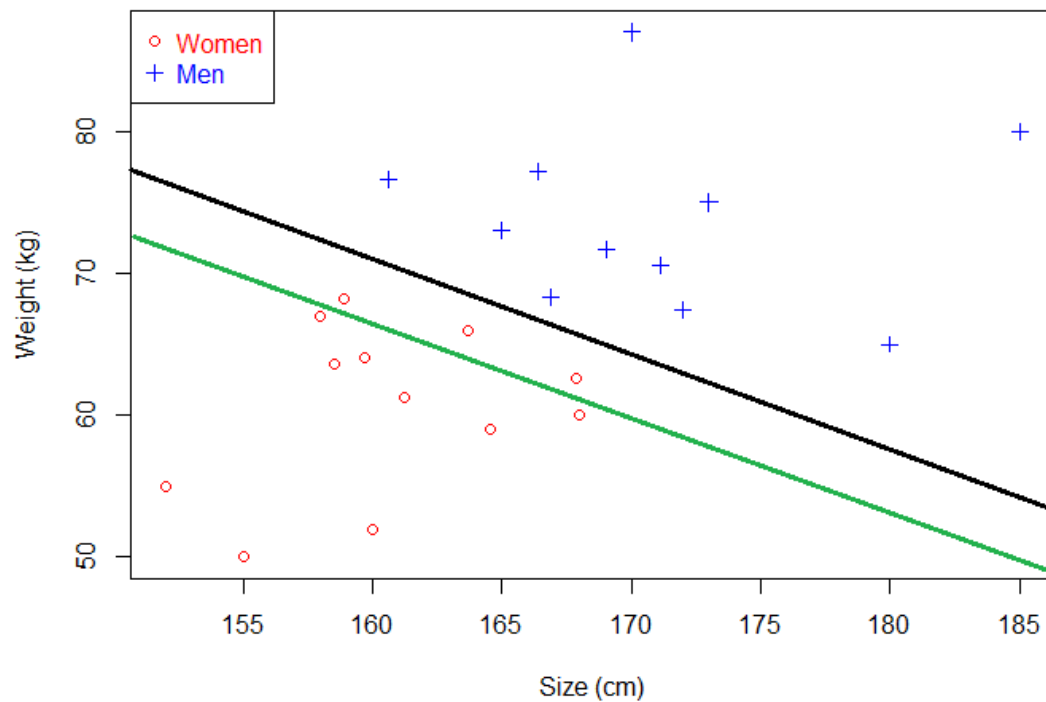


Figure 6: The black hyperplane classifies more accurately than the green one

SVMs compare two hyperplanes (which one is better for separate data) base on **the margin**. Basically the margin is a no man's land. There will never be any data point inside the margin.

We can make the following observations:

- If an hyperplane is very close to a data point, its margin will be small.

- The further an hyperplane is from a data point, the larger its margin will be.

This means that the optimal hyperplane will be the one with the biggest margin.

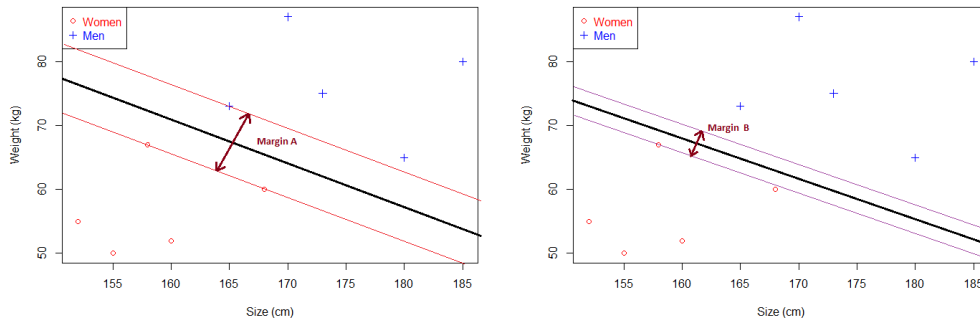


Figure 7: Margin B is smaller than Margin A, so Margin A would be the better one for separating data.

In practice, we can compute the distance between the hyperplane and the closest data point. Once we have this value, if we double it we will have the value of margin. After determining what the margin is and how to calculate it, we will solve the optimal problem (quadratic problem) to find the *support vectors* and compute the  $w$ ,  $b$  parameters for the optimal hyperplane from those *support vectors*.

We will not go on the mathematics behind SVMs, but you can read [this book](#).

## Pseudo-code

For hard margin classification, all we need to do is let mathematics do the job (solve the quadratic problem). The quadratic problem solver we used is available here [scikit-learn library](#). Below is an example code for using it.

---

```

1 import cvxopt.solvers
2 X, y = get_dataset(ls.get_training_examples)
3 m = X.shape[0]
```

```

4 # Gram matrix – The matrix of all possible inner products of X.
5 K = np.array([np.dot(X[i], X[j])
6   for j in range(m)
7   for i in range(m)]).reshape((m, m))
8 P = cvxopt.matrix(np.outer(y, y) * K)
9 q = cvxopt.matrix(-1 * np.ones(m))
10 # Equality constraints
11 A = cvxopt.matrix(y, (1, m))
12 b = cvxopt.matrix(0.0)
13 # Inequality constraints
14 G = cvxopt.matrix(np.diag(-1 * np.ones(m)))
15 h = cvxopt.matrix(np.zeros(m))
16 # Solve the problem
17 solution = cvxopt.solvers.qp(P, q, G, h, A, b)
18 # Lagrange multipliers
19 multipliers = np.ravel(solution['x'])
20 # Support vectors have positive multipliers.
21 has_positive_multiplier = multipliers > 1e-7
22 sv_multipliers = multipliers[has_positive_multiplier]
23 support_vectors = X[has_positive_multiplier]
24 support_vectors_y = y[has_positive_multiplier]

```

---

For computing w

---

```

1 def compute_w(multipliers, X, y):
2     return np.sum(multipliers[i] * y[i] * X[i]
3       for i in range(len(y)))

```

---

For computing b

---

```

1 def compute_b(w, X, y):
2     return np.sum([y[i] - np.dot(w, X[i])
3       for i in range(len(X))]) / len(X)

```

---

## 4 Experiments

### 4.1 Data preprocessing

We randomly split data into train (80) and test (20) subsets using `sklearn.model_selection.train_test_split`. The ground-truth labels of data will be get from the folder from original dataset and then it will be stored in **features/y\_train.pkl** and **features/y\_test.pkl** for using in python code.

### 4.2 Feature extraction

The BOV vocabulary is built up by using all images in dataset. We get about 40 - 50 key points and its descriptors in every image, then run K-Means with **k = 500** on this key points space to get the vocabulary with 500 visual words. The vocabulary is stored in **bin/vocabulary.pkl**

After having the vocabulary, we calculate the histogram of visual words for all images in both training set and test set. The result (feature vectors of training and test set) is stored in **features/X\_train.pkl** and **features/X\_test.pkl**

### 4.3 KNN test

We simply run KNN on test set with different k (from 1 to 100) and end up with result below.

### 4.4 SVMs test

We try to train SVM model with different parameter values, different kernels and aproachs for multi-classification. All models are saved in **bin/svm-models** with format "**multi-class\_kernel\_c\_degree\_gamma**". The value of parameter

is changing following the conclusion from theory we've learn in the class and our experiences in the testing process.

Here is some results (for more detail, you can read the log files from **log/svm**).

<i>method</i>	<i>C</i>	<i>training – time</i>	<i>training – accuracy</i>	<i>test – accuracy</i>
ovo	0.001	39.568	0.126	0.044
	0.01	36.864	0.326	0.151
	0.1	41.069	0.998	0.142
	0.2	44.101	1	0.142
	0.3	41.413	1	0.142
	100	40.766	1	0.142
	1000	44.251	1	0.142
ovr	0.001	2.713	0.406	0.060
	0.01	3.411	0.840	0.101
	0.1	5.355	0.977	0.104
	0.2	5.956	0.986	0.085
	0.3	6.241	0.988	0.079
	100	7.320	0.975	0.038
	1000	7.654	0.974	0.041

Table 1: SVM result with linear kernel.

<i>method</i>	<i>C</i>	<i>degree</i>	<i>training – time</i>	<i>training – accuracy</i>	<i>test – accuracy</i>
ovo	0.1	2	39.325	0.114	0.009
	0.1	3	34.676	0.115	0.009
	1	2	46.896	0.117	0.047
	1	3	51.521	0.116	0.035
	10	3	67.235	0.130	0.075
	100	2	40.766	0.998	0.157
	100	3	44.251	0.216	0.138
	1000	2	75.657	1.000	0.154
	1000	3	72.126	0.998	0.126

Table 2: SVM result with poly kernel.



<i>method</i>	<i>C</i>	<i>gamma</i>	<i>training – time</i>	<i>training – accuracy</i>	<i>test – accuracy</i>
ovo	0.1	2	39.325	0.114	0.009
	0.1	3	34.676	0.115	0.009
	1	2	46.896	0.117	0.047
	1	3	51.521	0.116	0.035
	10	3	67.235	0.130	0.075
	100	2	40.766	0.998	0.157
	100	3	44.251	0.216	0.138
	1000	2	75.657	1.000	0.154
	1000	3	72.126	0.998	0.126

Table 3: SVM result with rbf kernel.

### Some theories have been proven

- "One Vs One" accuracy is higher than "One Vs Rest" (shown in ??)
- "One Vs Rest" accuracy is faster than "One Vs One" (shown in ??)
- If the value of parameter C for SVM soft-margin is great enough, soft margin will behave like hard margin (shown in ?? the One Vs One method,  $C = 0.2$  or greater we have the same result (and it also is the hard-margin result)).
- The more complex the model is, the more time takes to train and the trained model can be cause of over-fitting (the data is more simpler than what the model represent). (In ?? when *degree* becomes greater or in ?? when *gamma* becomes greater, it takes a lot of time to train and the trained model fits almost 100% training data)

## 5 Conclusions

The result accuracy is quite really small and the reason we guess is the simple feature extractor (BOV) we have been implemented. The main goal of the project (summary what we've learned, work on a real dataset to see what we really learned

and how those algorithms work together) was archived. Next, we will focus on changing our feature extractor to improve the result.

## 6 References

- [1] F. Liu, M. Gleicher, H. Jin and A. Agarwala, ‘Content-preserving warps for 3d video stabilization’, *ACM Transactions on Graphics*, 2009.
- [2] C. Demant, B. Streicher-Abel and P. Waszkewitz, *Industrial image processing: Visual quality control in manufacturing*. Springer, 2013.
- [3] H. G. Jung, D. S. Kim, P. J. Yoon and J. Kim, *Structural, Syntactic and Statistical Pattern Recognition*. Springer Berlin / Heidelberg, 2006.
- [4] Wikipedia, *Support vector machine*, 2004. [Online]. Available: [https://en.wikipedia.org/wiki/Support\\_vector\\_machine](https://en.wikipedia.org/wiki/Support_vector_machine).

## 7 Appendices