

# Introduction to Artificial Intelligence: Assignment #3

Autumn 2015

Due: Nov 6rd, 23:59:59 CST (UTC +8).

## 1. K-Nearest Neighbor

In this problem, we will play with K-Nearest Neighbor (KNN) algorithm and try it on real-world data. Implement KNN algorithm (in *knn.m*), then answer the following questions.

- (a) In *knn\_exp.m*, try KNN with different K (you should at least experiment  $K = 1, 10$  and 100) and plot the decision boundary.

You are encouraged to vectorize<sup>1</sup> your code, otherwise the experiment time might be extremely long. You may find the MATLAB build-in functions *pdist2*, *sort*, *max* and *hist* useful. Also, you can use the function *eudist2*<sup>2</sup> written by Prof. Deng Cai<sup>3</sup>.

- (b) We have seen the effects of different choices of K. How can you choose a proper K when dealing with real-world data ?
- (c) Now let us use KNN algorithm to hack the CAPTCHA of a website<sup>4</sup> that we are all familiar with:



Finish *hack.m* to recognize the CAPTCHA image using KNN algorithm.

You should label some training data yourself, and store the training data in *hack\_data.mat*. Helper functions *extract\_image* and *show\_image* are give for your convenience.

Remember to submit *hack\_data.mat* along with your code and report.

<sup>1</sup>[http://www.mathworks.cn/cn/help/matlab/matlab\\_prog/vectorization.html](http://www.mathworks.cn/cn/help/matlab/matlab_prog/vectorization.html)

<sup>2</sup><http://www.cad.zju.edu.cn/home/dengcai/Data/code/EuDist2.m>

<sup>3</sup>Prof. Deng Cai is an expert on MATLAB, you can find all his code at <http://www.cad.zju.edu.cn/home/dengcai/Data/data.html>. You can learn how to write fast MATLAB code by reading his code.

<sup>4</sup><http://jwbinfosys.zju.edu.cn/default2.asp>

## 2. Decision Tree and ID3

Consider the scholarship evaluation problem: selecting scholarship recipients based on gender and GPA. Given the following training data:

Gender	GPA	Scholarship	Count
F	Low	+	10
F	High	+	95
M	Low	+	5
M	High	+	90
F	Low	-	80
F	High	-	20
M	Low	-	120
M	High	-	30

Draw the decision tree that would be learned by ID3 algorithm and annotate each non-leaf node in the tree with the information gain attained by the respective split.

## 3. K-Means Clustering

Finally, we will run our first unsupervised algorithm – k-means clustering. Implement k-means algorithm (in *kmeans.m*), then answer the following questions.

Note that there are different kind of methods to setup initial cluster centers for k-means algorithm, we will use a simple one – randomly choose  $K$  samples from dataset as initial cluster centers.

- Run your k-means algorithm on *kmeans\_data.mat* with the number of clusters  $K$  set to 2. Repeat the experiment 1000 times. Use *kmeans\_plot.m* to visualize the process of k-means algorithm for the two trials with largest and smallest SD (sum of distances from each point to its respective centroid).
- You should observe the issue that the outcome of k-means algorithm is very sensitive to cluster centroids initialization from the above experiment. How can we get a stable result using k-means?
- Run your k-means algorithm on the digit dataset *digit\_data.mat* with the number of clusters  $K$  set to 10, 20 and 50. Visualize the centroids using *show\_digit.m*. You should be able to observe that k-means algorithm can discover the patterns in dataset without any label information.
- Another important application of k-means is Vector quantization<sup>5</sup>. Vector quantization is a classical quantization technique from signal processing. It works by dividing a large set of points (vectors) into groups, then representing the data points by their group centroid points, as in k-means and some other clustering algorithms.

---

<sup>5</sup>[https://en.wikipedia.org/wiki/Vector\\_quantization](https://en.wikipedia.org/wiki/Vector_quantization)

Here we will use vector quantization to do image compression. By clustering image pixel value into  $K$  groups, we can represent each pixel with  $\log(K)$  bits, instead of 24 bits (RGB, each channel has 8bit depth).

Finish *vq.m*. Compress images with  $K$  set to 8, 16, 32 and 64. I have provided you some sample images, however use your own photos is encouraged.

What is the compress ratio if we set  $K$  to 64 (Optionally, you can compute the compress ratio using Huffman encoding) ?

## 4. Spectral Clustering

In this problem, we will try a dimensionality reduction based clustering algorithm – Spectral Clustering. Implement Spectral Clustering algorithm in *spectral.m*, then answer the following questions.

- (a) We will first experiment Spectral Clustering on synthesis data (in *spectral\_exp1.m*). Finish *knn\_graph.m* to construct the KNN graph<sup>6</sup>  $W$ , then test your algorithm on data *cluster\_data.mat*. Visualize the clustering result and compare it with k-means.

Note that we only care about the local connectivity of data points, therefore edges between far away points should be discarded. You should take care with this issue when implementing *knn\_graph.m* and choose a proper threshold.

- (b) Now let us try Spectral Clustering on real-world data (in *spectral\_exp2.m*). The TDT2 corpus consists of data collected during the first half of 1998 and taken from 6 sources, including 2 newspapers (APW, NYT), 2 radio programs (VOA, PRI) and 2 television programs (CNN, ABC).

Try Spectral Clustering on a subset of TDT2 corpus (in *TDT2\_data.mat*) and compare the result with k-means. You should evaluate the quality of clustering using accuracy and normalized mutual information<sup>7</sup>.

For this part of problem, use *constructW.m* to construct the graph  $W$ , and choose proper parameters. As usual, you should preprocessing the data and repeat the experiments multiple times then take the average.

---

Please submit your homework report in **pdf** format, with all your code in a zip archive. Your files are suggested to be named as studentnum.name.pdf/zip/rar.

---

<sup>6</sup>Connect each data point with its  $K$  nearest neighbors, you can use binary edge weights or compute edge weights using Heat kernel.

<sup>7</sup>See <http://www.cad.zju.edu.cn/home/dengcai/Data/Clustering.html> and Deng Cai, Xiaofei He, and Jiawei Han, "Document Clustering Using Locality Preserving Indexing", in IEEE TKDE, 2005, for more details.