

# CS 383 - Machine Learning

## Assignment 2 - Clustering Winter 2017

### Introduction

In this assignment you'll work on clustering data.

You may not use any function from the Matlab ML library in your code. Look at the *Matlab Functions* section on Blackboard for a list of functions that are ok to use.

In particular for this assignment you **MAY NOT** use Matlab functions like:

- k-means functions

But you **MAY** use basic statistical functions like:

- std
- mean
- cov
- eig

*As a reminder, make sure to clear out old variables prior to running your script.*

### Grading

Although all assignments will be weighed equally in computing your homework grade, below is the grading rubric we will use for this assignment:

Part 3 (k-Means) Report	30pts 5pts
<b>TOTAL</b>	<b>35pts</b>

Table 1: Grading Rubric

# DataSets

**Pima Indians Diabetes Data Set** In this dataset of 768 instances of testing Pima Indians for diabetes each row has the following information (1 class label, 8 features).

1. Class Label (-1=negative,+1=positive)
2. Number of times pregnant
3. Plasma glucose concentration
4. Diastolic blood pressure (mm Hg)
5. Triceps skin fold thickness (mm)
6. Insulin ( $\mu$ U/ml)
7. Body mass index ( $kg/m^2$ )
8. Diabetes pedigree function
9. Age (yrs)

Data obtained from: <https://archive.ics.uci.edu/ml/datasets/Pima+Indians+Diabetes>

# 1 Clustering

Next we're going to cluster this same data using k-means!

**Write a script that:**

1. Reads in the data
2. Standardizes the data
3. Performs k-means clustering **using just the 6th and 7th feature of the data (BMI and Pedigree, respectively) with  $k=2$**  and using the Euclidean (L2) distance.

In addition, in your k-means code you'll want to visualize the progress of the algorithm (this will be part of your report):

1. Plot the initial setup
  - (a) Data points are red 'x'
  - (b) Cluster centers are blue 'o'
2. Plot the initial cluster assignments
  - (a) Cluster 1 = red
  - (b) Cluster 2 = blue
  - (c) Data points are as 'x' (according to their assigned color)
  - (d) Cluster centers are as 'o' (according to their assigned color)
3. Plot the final cluster assignments
  - (a) Cluster 1 = red
  - (b) Cluster 2 = blue
  - (c) Data points are as 'x' (according to their assigned color)
  - (d) Cluster centers are as 'o' (according to their assigned color)
  - (e) Title should indicate how many iterations it took to get there

Your figures should end up similar to Figures 1-3.

## Implementation Details

1. Seed the random number generator with zero (do this right before running your k-means). You can use Matlab's **rng** function to do this.
2. Randomly select two data instances and use them for the initial seeds (since we'll do  $k = 2$ ). I suggest you use `randperm` to randomize the indices and use the first two to select the observations you will use to initialize your reference vectors/cluster centers.
3. Use the L2 distance to measure the distance between observations and reference vectors.

4. Terminate the EM process when the sum of magnitude of change of the cluster centers (from the previous iteration to the current one) is less than  $eps$  (which is a Matlab defined variable related to the possible floating-point precision). That is, when  $\sum_{i=1}^k d(a_i(t-1), a_i(t)) < \epsilon$  where  $k$  is the number of clusters,  $a_i(t)$  is the reference vector for cluster  $i$  at time  $t$  and  $d(x, y)$  is the L1 distance between vectors  $x$  and  $y$  (as defined in the *Similarity and Distance Functions* link on BBlearn).
5. Write your code in such a way that it could work for any value of positive integer  $k$ , and any number of features,  $D$ . However you only have to plot the first two features and only have to plot *two* clusters.

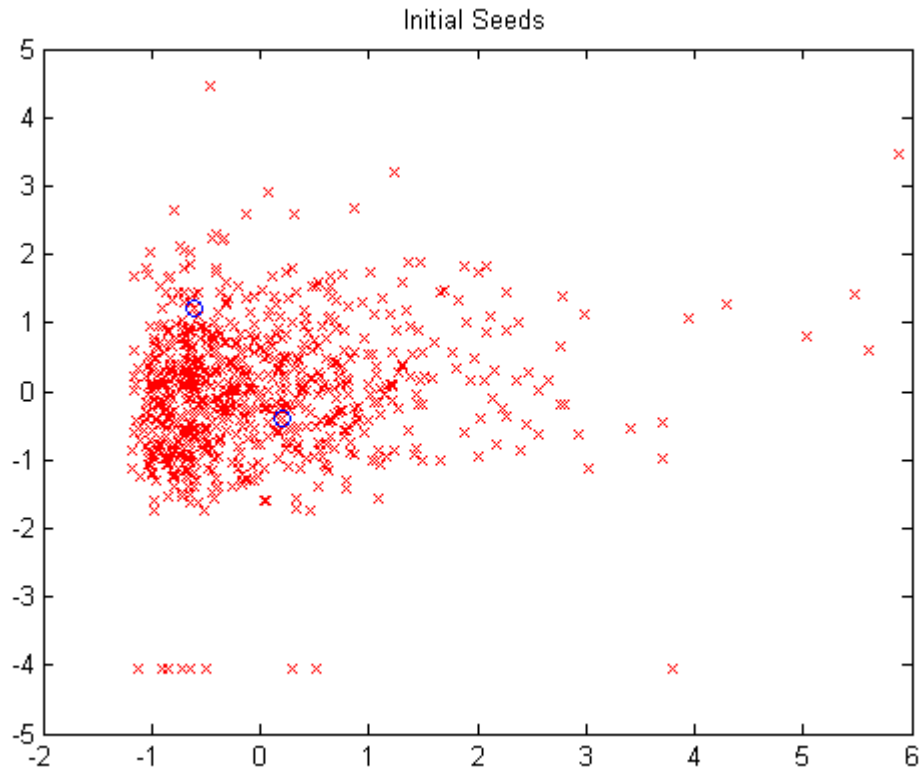


Figure 1: Seeds

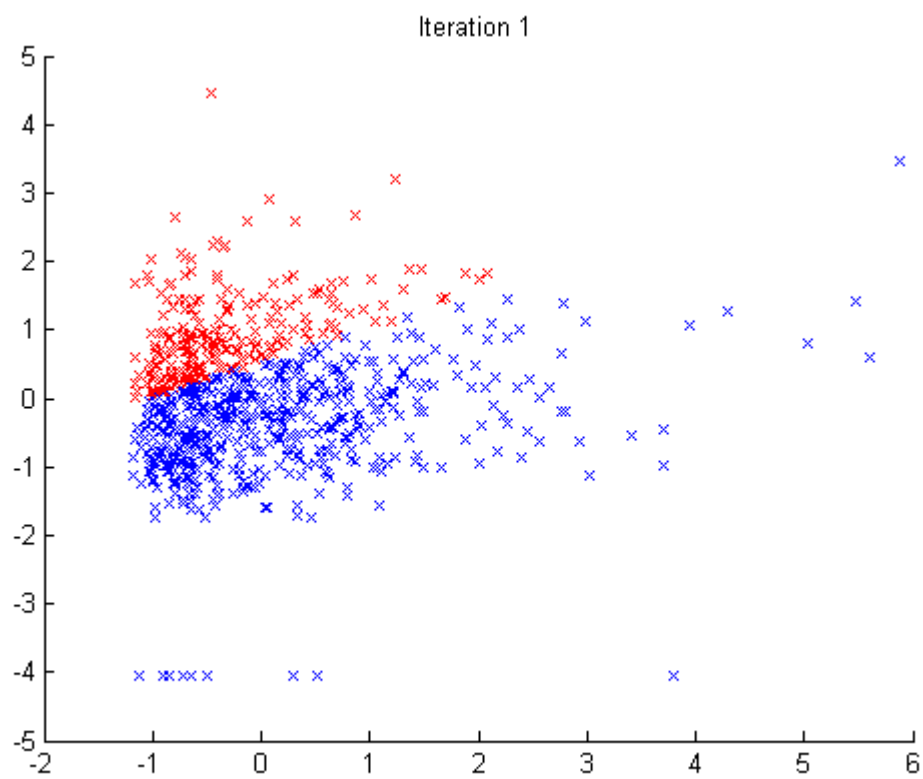


Figure 2: Initial Clustering

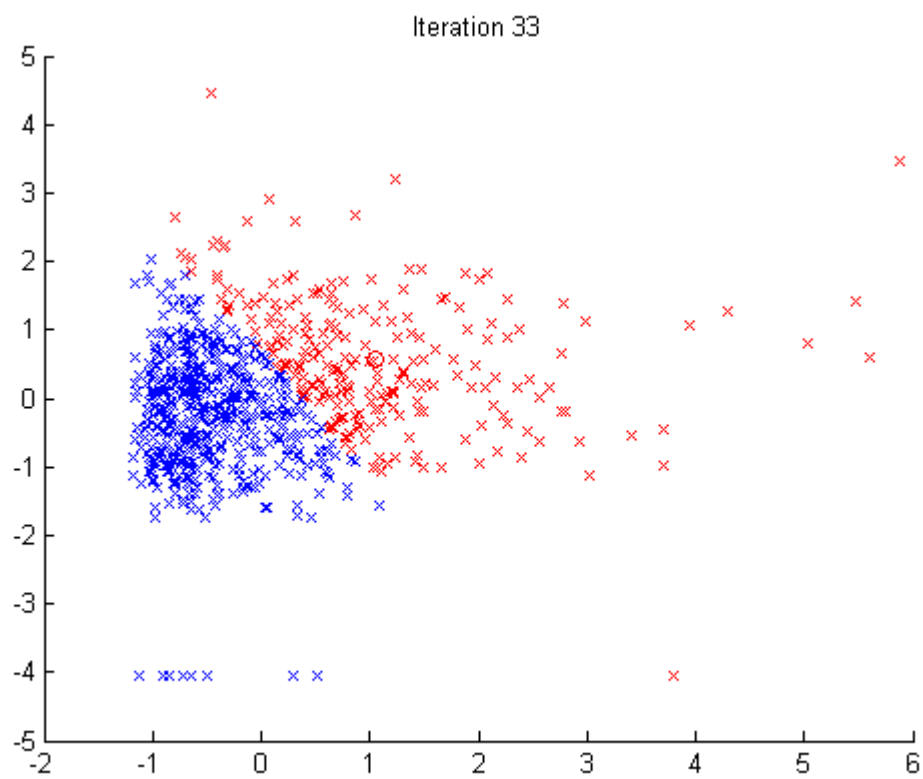


Figure 3: Final Clustering after 33 iterations

# Submission

For your submission, upload to Blackboard a single zip file containing:

1. PDF Writeup
2. Source Code
3. readme.txt file

The readme.txt file should contain information on how to run your code to reproduce results for each part of the assignment.

The PDF document should contain the following:

1. Part 1: The visualization of the k-means clustering process including:
  - (a) The initial setup visualization
  - (b) The initial cluster assignment visualization
  - (c) The final cluster assignment visualization

and report how many iterations it took for your algorithm to terminate.