

Lab 6

Albert Segarra Roca (S3255050), Carlos Humberto Paz Rodríguez (S3040577)

Group 16

Pattern Recognition

FMNS • RUG

October 26, 2016

1 Assignment 1

1.1

As we can see in the code provided in section 3.1 the implementation is straightforward. First, we initialize the prototypes with random points from the data, and we assign each point in data to a prototype by taking that with minimum distance to the point. Then, we start the loop by computing the average point among those of the same prototype, and we reassign all points to according to the new prototypes. This is looped until there is no difference between the previous assignment and the current one, meaning no point has been reassigned.

Note that we store and return a `prototypesHistory` which contains the list of prototype points for every iteration. This is later used to plot the arrows and intermediate positions of the prototypes.

The implementation for the plots can be found in section 3.2. For the data points and final prototypes it is straightforward and its results for $k = 2, 4, 8$ can be seen in Figures 1,2,3 respectively. Note that different colors indicate different clusters.

As we can see, for $K = 2$ the points are divided in two clear clusters, the one in the left and the one in the right. For $k = 4$ and $k = 6$, these two clusters are subdivided in 2 and 3 clusters respectively.

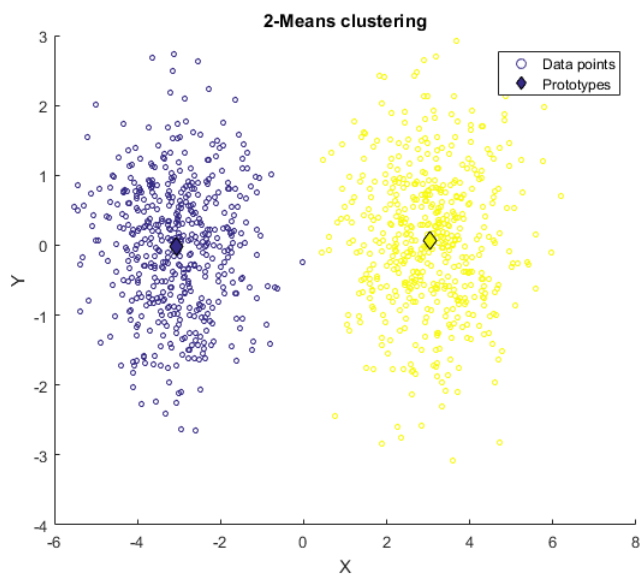


Figure 1: 2-means clustering data points and prototypes

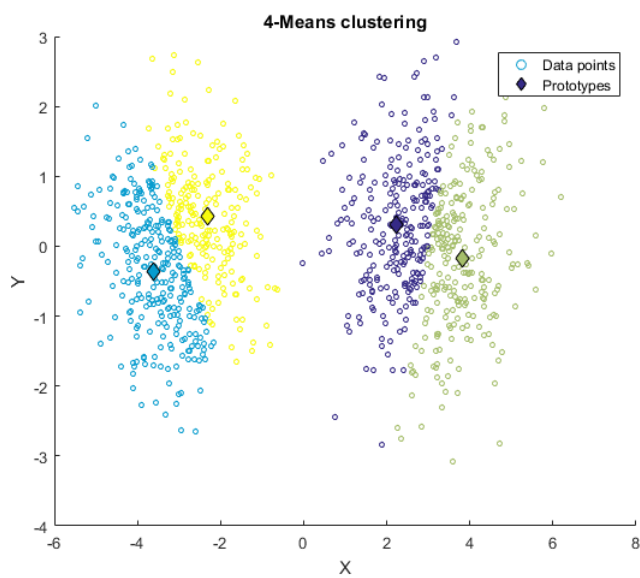


Figure 2: 4-means clustering data points and prototypes

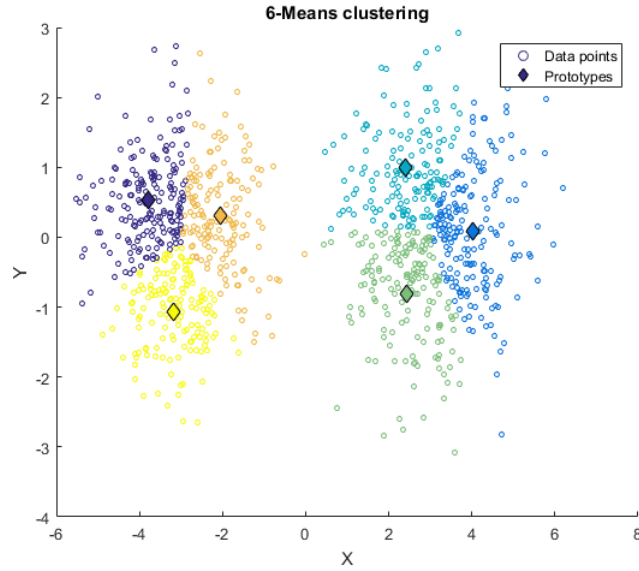


Figure 3: 6-means clustering data points and prototypes

In order to plot the prototype movements, we use the `prototypesHistory` returned by our `kmeans` function. We plot the initial, intermediate and final prototypes separately and with different markers to be able to differentiate them. For the intermediate prototypes, we loop through the `prototypesHistory` array and we print both a mark and an arrow from the previous prototype to the current prototype. We do so by using the `plotArrows` function, which plots an arrow from a prototype to another. Again, its results can be seen in Figures 4 ($k = 2$), 5 ($k = 4$) and 6 ($k = 6$).

We can see in this case how the prototypes quickly move to a much better prototype in the first step and then take smaller steps each time as the algorithm advances. We can also note in general that it takes more iterations for the algorithm to converge for bigger K .

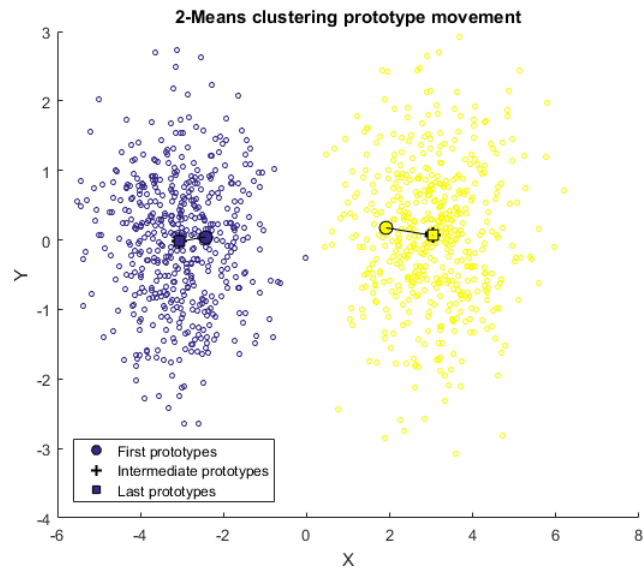


Figure 4: 2-means clustering prototype movement

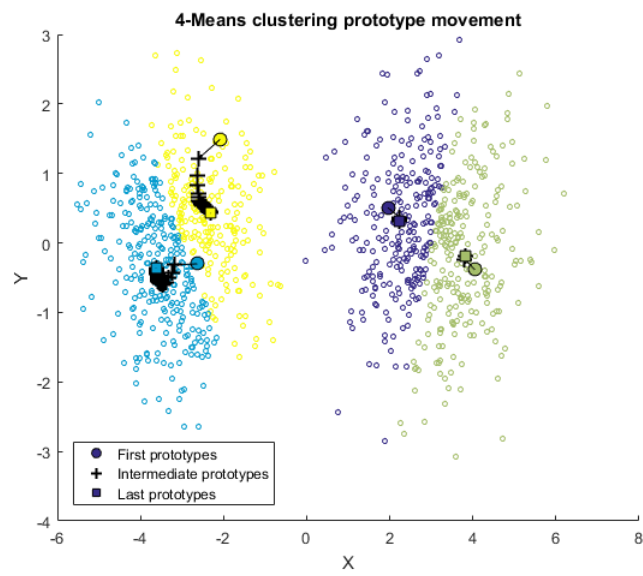


Figure 5: 4-means clustering prototype movement

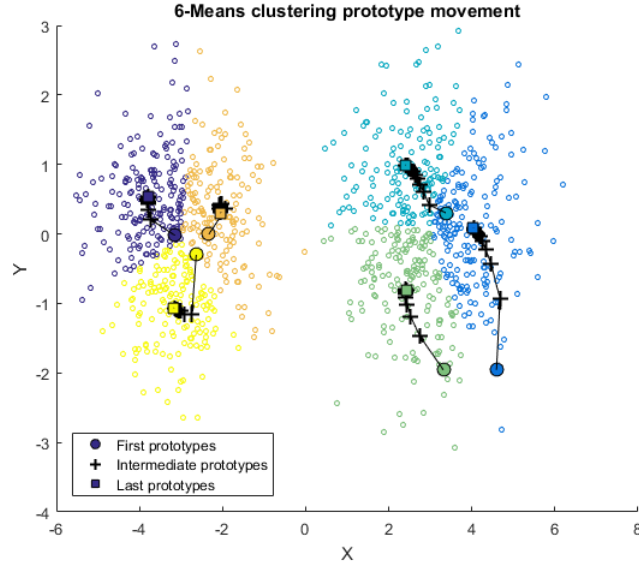


Figure 6: 6-means clustering prototype movement

1.2

In order to compute $J(k)$, the quantization error function, we need to compute half the sum of squares of euclidean distances between each point and its prototype. Then we compute $R(k)$ and $D(k)$ with the straightforward formulas indicated in the statement. As we can see in the code provided in section 3.3, we have chosen k_{max} to be 10, which we found to be a reasonably high value for the number of clusters such that higher values would not result in better clustering. In order to have more precise results, we compute each $J(k)$ as the average of the quantization error for 10 repetitions of the **k-means** algorithm, as the results can vary depending on the initialization of the prototypes.

- a) As we can see in the plot of $D(k)$ (Figure 7) the function has a maximum for $k = 2$, and does roughly only decrease for higher values of k . We can conclude then that the optimal value for the number of clusters for this dataset is $k_{opt} = 2$. This makes sense if we look at figure 1, as there is a clear distinction between points belonging to the left cluster and the right cluster. Adding more clusters does not make many sense in this dataset, because it does only partition this two main clusters in smaller pieces.

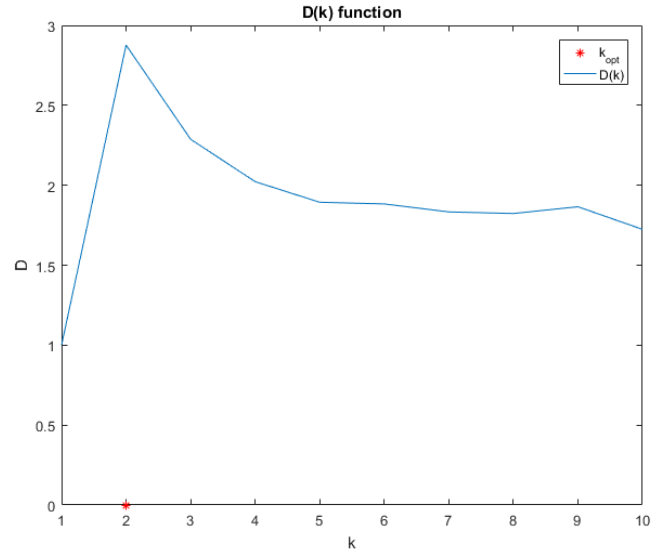


Figure 7: $D(k)$ function with k_{opt}

- b) Also when we plot $J(k)$ and $R(k)$ (Figure 8) we see how k_{opt} is achieved when both functions have the biggest distance and $J(k)$ is below $R(k)$, thus $D(k) = \frac{R(k)}{J(k)}$ is optimized because $R(k)$ has the highest value as compared to $J(k)$.

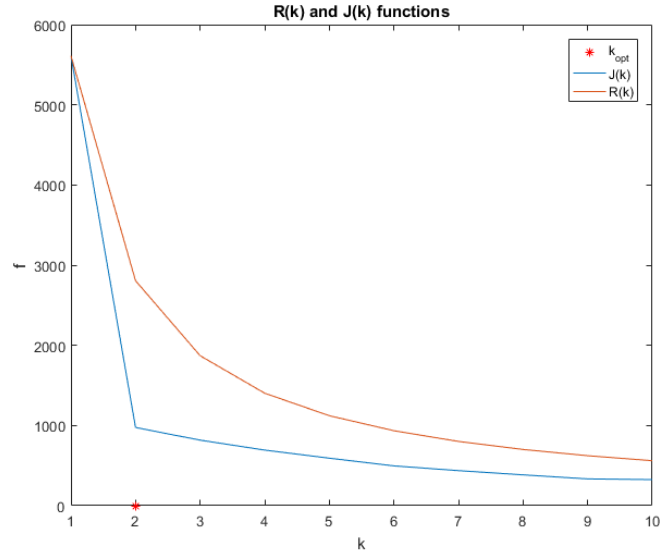


Figure 8: $R(k)$ and $J(k)$ functions with k_{opt}

1.3

- a) In order to implement the **k-means++** initialization, we use the algorithm described in the statement. The only interesting non-trivial step is how to pick a data point as a new prototype with probability proportional to the distance to its closest prototype. As seen in function **selectPrototypeDx** in section 3.1 (lines 14-32), to accomplish that we first pick a random number R in the range $(0, \text{sum}(D))$. This random number is then used to pick a point X in the cumulative function D_c of D such that this $D_c(X) \geq R$. Because each point X takes a span of the range of the random number which is proportional to $D(X)$, the point has a probability to be selected which is proportional to its span, thus proportional to $D(X)$.
- b) To compare the two algorithms, we run them both 20 times and compute the mean and standard deviations of their quantization errors. The implementation can be found in section 3.4.
- (a) **k-means quantization error:**
Mean: **8.3610**
Standard deviation: **0.6889**
- (b) **k-means++ quantization error:**
Mean: **6.7081**
Standard deviation: **0.4615**

Now as we can see, the **k-means++** algorithm clearly has a lower mean quantization error than the **k-means** algorithm, although the difference is not very big. This is because **k-means++** uses a better initialization for the prototypes, which is actually very important for the algorithm because it tends to get stuck in local optima.

- c) Now, to compute the desired p-value and null hypothesis test, we use Matlab's **ttest2** function, for a two-sample t-student test, without the assumption of equal variances (Welch's t-test). The resulting p-value is with the null hypothesis that the two data vectors are from populations with equal mean.

p-value: 2.5151e-10

As we can see, this value is very low, which indicates that the probability that they have equal means is very low. Also, with the first return value of **ttest2**, which returns 1 in our case, we can reject the null hypothesis (equal means) with a significance of 5% that and thus with that significance we can say that the means are different. Because **k-means++** quantization error mean is lower we can conclude that **k-means++** outperforms **k-means** algorithm.

2 Assignment 2

2.1

The full implementation of the batch neural gas algorithm can be found in section 3.5. The only non-trivial step is explained in the comments, which is the computation of the ranking for each prototype, for which sorting is used to compute the number of prototypes with smaller distance to the point.

2.2

In order to make all the plots, we used the function provided in section 3.6.

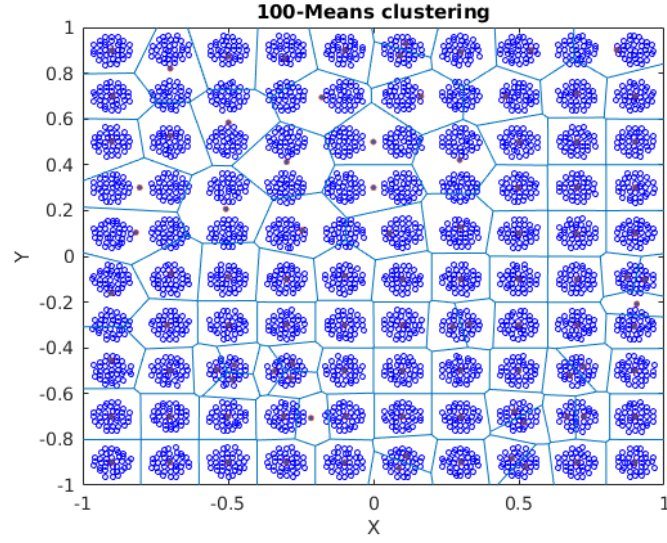


Figure 9: 100-means clustering for checkerboard data

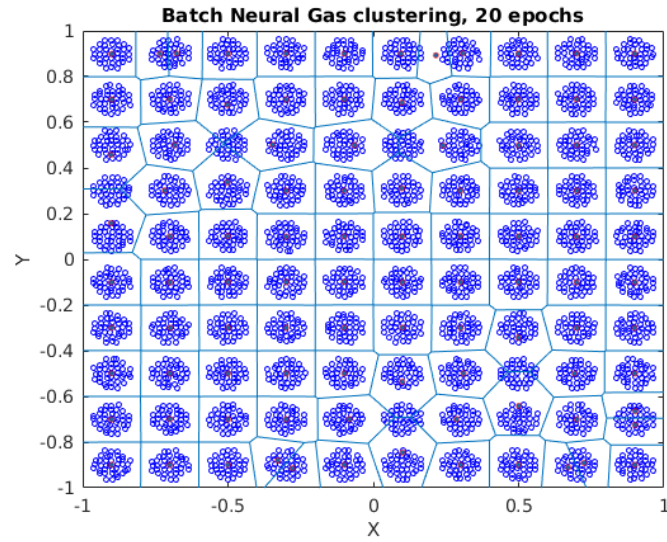


Figure 10: Batch Neural Gas clustering with 20 epochs for checkerboard data

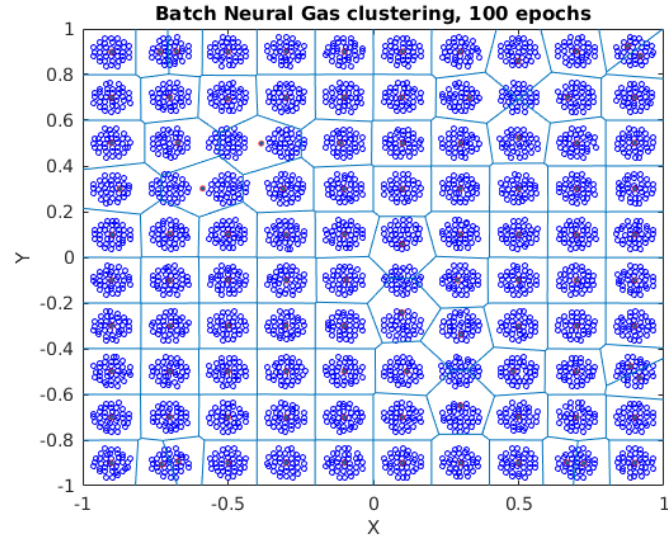


Figure 11: Batch Neural Gas clustering with 100 epochs for checkerboard data

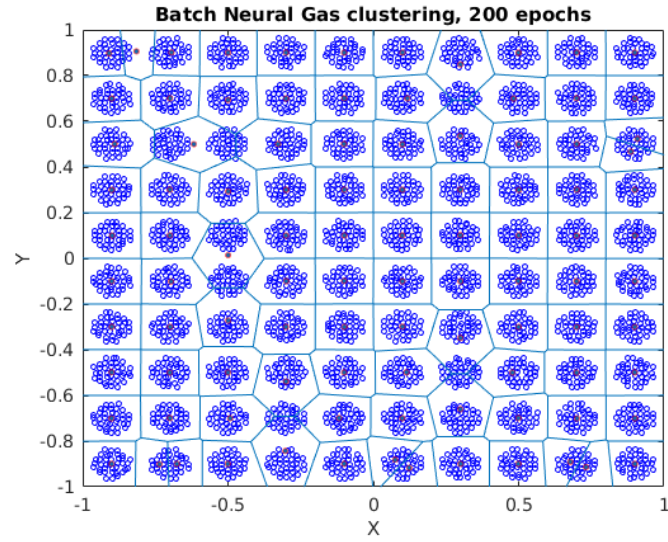


Figure 12: Batch Neural Gas clustering with 200 epochs for checkerboard data

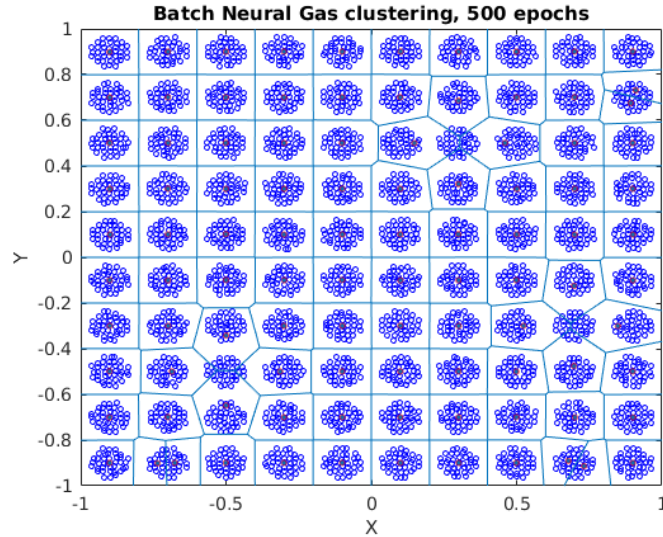


Figure 13: Batch Neural Gas clustering with 500 epochs for checkerboard data

2.3

As we can see by comparing figures 9 and 10, even with only 20 epochs the batch neural gas algorithm already outputs an apparently better clustering than the k-means algorithm. The results get even better as we increase the number of epochs for the BNG algorithm, as we can see by comparing figures 10, 11, 12 and 13, although it doesn't improve very significantly.

The difference between both algorithms is that **k-means** tries to minimize directly the quantization error, whereas **BNG** minimizes a cost function instead with a gradient descent.

The improvement seen in the BNG algorithm as compared to the **k-means** algorithm is because the latter usually gets stuck in local optima, and it is difficult for it to find the global optimum. Also, it is much more sensitive to the initialization of the prototypes, as compared to the BNG algorithm, and also it is very sensitive to outliers and the variation of size and density of the clusters, which makes it not adequate for some data sets. Both algorithms share the disadvantage that we need to give it the number of desired clusters, which is sometimes difficult to guess beforehand.

So, generally, the BNG algorithm gives better results, although **k-means** algorithm works quite well for gaussian-shaped clusters, but still the BNG algorithm is much more computationally expensive than the **k-means** algorithm, which is generally quite fast.

3 Appendix

3.1 k-means function implementation (with k-means++)

```

1 function [classes, prototypesHistory] = kmeans(data, k, is_kmeanspp, initialize_file)
2     if nargin < 4
3         initialize_file = false;
4         if nargin < 3
5             is_kmeanspp = false;
6         end
7     end
8
9     function classes = reassign(data, prototypes)
10        % Compute distances between every point and prototype
11        distances = pdist2(data, prototypes);
12        % For every point, get the index of the prototype with min distance
13        [~, classes] = min(distances, [], 2);
14    end
15
16    function prototype = selectPrototypeDx(Dsq)
17        % Selects a prototype x with probability proportional to
18        % D(x, closestPrototype(x))
19        sumDistances = sum(Dsq);
20        randX = rand*sumDistances; % Get random in range (0, sumDistances)
21        i = 1;
22        accSum = Dsq(i);
23
24        while randX > accSum
25            i = i + 1;
26            accSum = accSum + Dsq(i);
27        end
28
29        % Pick prototype in the position X where the cumulative sum is
30        % greater than the random. Obviously elements with greater D(x)
31        % will span a bigger portion of the range of the random number,
32        % and thus the probability of selecting them will be higher
33        prototype = data(i, :);
34    end
35
36    n = size(data, 1);
37
38    prototypesHistory = [];
39
40    if is_kmeanspp
41        % k-means++ prototype initialization
42        prototypes = zeros(k, 2);
43        prototypes(1, :) = data(randi(n), :);
44        diff = data - repmat(prototypes(1, :), n, 1);
45        % Not using pdist^2 cause it would be redundant and inefficient
46        Dsq = diff(:, 1).^2 + diff(:, 2).^2;
47
48        for prototypeIndex = 2:k-1
49            prototypes(prototypeIndex, :) = selectPrototypeDx(Dsq);
50            classes = reassign(data, prototypes(1:prototypeIndex, :));
51            diff = prototypes(classes, :) - data;
52            Dsq = diff(:, 1).^2 + diff(:, 2).^2;
53        end
54
55        if k > 1
56            prototypes(k, :) = selectPrototypeDx(Dsq);
57        end
58    elseif initialize_file
59        sbrace = @(x,y)(x{y});
60        fromfile = @(x)(sbrace(struct2cell(load(x)),1));
61        prototypes=fromfile('clusterCentroids.mat');

```

```

62 else
63     % Compute initial prototypes selecting k random points from data
64     prototypes = data(randi(n, k, 1), :);
65
66 end
67
68 % Assign points to the prototype with minimum distance
69 classes = reassign(data, prototypes);
70
71 someReassigned = true;
72 it = 1;
73 while someReassigned
74     prototypesHistory(:, :, it) = prototypes;
75     it = it + 1;
76
77     prototypes = zeros(k, 2); % Assume 2D data
78     for i = 1:n % Compute sum vector of all vectors belonging to a specific prototype
79         prototypes(classes(i), :) = prototypes(classes(i), :) + data(i, :);
80     end
81
82     % Compute mean vector dividing the sum vector by the number of points in
83     % the prototype
84     prototypes = prototypes ./ repmat(hist(classes, 1:k)', 1, 2);
85
86     auxClasses = classes;
87     classes = reassign(data, prototypes);
88     someReassigned = ~isequal(auxClasses, classes);
89 end
90 prototypesHistory(:, :, it) = prototypes;
91
92 end

```

3.2 Plotting functions for exercise 1.1

```

1 function plotArrows(previousPrototypes, prototypes)
2     for i = 1:k
3         plot_arrow(previousPrototypes(i, 1), previousPrototypes(i, 2), prototypes(i, 1),
4                     prototypes(i, 2));
5     end
6 end

```

```

1 function kmeansplotter(K)
2     data = load('kmeans1');
3     data = data.kmeans1;
4
5     [classes, prototypesHistory] = kmeans(data, K);
6     iterations = size(prototypesHistory, 3);
7
8     % Plot 1: points and prototypes
9     figure(1);
10    scatter(data(:, 1), data(:, 2), 10, classes, 'DisplayName', 'Data points');
11    hold on;
12    lastPrototypes = prototypesHistory(:, :, iterations);
13    scatter(lastPrototypes(:, 1), lastPrototypes(:, 2), ...
14            70, [1:K]', 'filled', 'd', 'MarkerEdgeColor', [0 0 0], 'DisplayName', 'Prototypes');
15    xlabel('X');
16    ylabel('Y');
17    tt = sprintf('%d-Means clustering', K);
18    title(tt);
19    legend('show');

```

```

20 hold off;
21
22 % Plot 2: prototype movement
23 figure(2);
24 scatter(data(:, 1), data(:, 2), 10, classes);
25 hold on;
26 firstPrototypes = prototypesHistory(:, :, 1);
27 firstScatter = scatter(firstPrototypes(:, 1), firstPrototypes(:, 2), ...
28     70, [1:K]', 'filled', 'o', 'MarkerEdgeColor', [0 0 0]);
29
30 for i = 2:iterations-1
31     previousPrototypes = prototypesHistory(:, :, i-1);
32     prototypes = prototypesHistory(:, :, i);
33
34     intermediateScatter = scatter(prototypes(:, 1), prototypes(:, 2), ...
35         70, '+', 'MarkerEdgeColor', [0 0 0], 'LineWidth', 1.5);
36     plotArrows(previousPrototypes, prototypes);
37 end
38
39 lastScatter = scatter(lastPrototypes(:, 1), lastPrototypes(:, 2), ...
40     70, [1:K]', 'filled', 's', 'MarkerEdgeColor', [0 0 0]);
41 plotArrows(prototypesHistory(:, :, iterations-1), lastPrototypes);
42
43 xlabel('X');
44 ylabel('Y');
45 tt = sprintf('%d-Means clustering prototype movement', K);
46 title(tt);
47
48 if K == 1
49     legend([firstScatter lastScatter], ...
50         'First prototypes', 'Last prototypes', 'Location', 'southwest');
51 else
52     legend([firstScatter intermediateScatter lastScatter], ...
53         'First prototypes', 'Intermediate prototypes', 'Last prototypes', 'Location', '
54         southwest');
55 end
56 hold off;

```

3.3 Function for quantization error analysis and plotting

```

1 function qe = quantization_error(data, classes, prototypes)
2     % I do not use pdist because it would be redundant and inefficient to
3     % square a sqrt
4     qe = sum((data(:, 1) - prototypes(classes, 1)).^2 + ...
5         (data(:, 2) - prototypes(classes, 2)).^2)/2;
6 end

```

```

1 function kmeansquantization()
2 REPETITIONS = 30;
3
4 data = load('kmeans1');
5 data = data.kmeans1;
6 n = size(data, 1);
7
8 k_max = min(10, n);
9
10 J = zeros(k_max, 1);
11 for k = 1:k_max
12     for time = 1:REPETITIONS
13         [classes, prototypesHistory] = kmeans(data, k);

```

```

14         lastPrototypes = prototypesHistory(:, :, size(prototypesHistory, 3));
15         J(k) = J(k) + quantization_error(data, classes, lastPrototypes);
16     end
17 end
18
19 J = J ./ REPETITIONS;
20
21 R = J(1).*(1./[1:k]'); % J(1)*k^(-2/d), d=2
22
23 D = R./J;
24
25 [~, k_opt] = max(D);
26
27 % Plot 1, D function with k_opt
28 figure(1);
29 d_plot = plot(D);
30 hold on;
31 k_opt_plot = plot(k_opt, 0, 'r*');
32 title('D(k) function');
33 xlabel('k');
34 ylabel('D');
35 xlim([1 10]);
36 legend([k_opt_plot d_plot], 'k_{opt}', 'D(k)');
37 hold off;
38
39 % Plot 2, R and J function with k_opt
40 figure(2);
41 j_plot = plot(J);
42 hold on;
43 r_plot = plot(R);
44 k_opt_plot = plot(k_opt, 0, 'r*');
45 title('R(k) and J(k) functions');
46 xlabel('k');
47 ylabel('f');
48 xlim([1 10]);
49 legend([k_opt_plot j_plot r_plot], 'k_{opt}', 'J(k)', 'R(k)');
50 hold off;
51 end

```

3.4 Script to compare k-means with k-means++

```

1 function qe = quantization_error(data, classes, prototypes)
2     % I do not use pdist because it would be redundant and inefficient to
3     % square a sqrt
4     qe = sum((data(:, 1) - prototypes(classes, 1)).^2 + ...
5              (data(:, 2) - prototypes(classes, 2)).^2)/2;
6 end
7
8 K = 100;
9 RUNS = 20;
10
11 load checkerboard.mat
12
13 data = checkerboard;
14
15 J20 = zeros(RUNS, 1);
16 J20pp = zeros(RUNS, 1);
17 for i = 1:RUNS
18     [classes, prototypesHistory] = kmeans(data, K, false);
19     lastPrototypes = prototypesHistory(:, :, size(prototypesHistory, 3));

```

```

13     J20(i) = quantization_error(data, classes, lastPrototypes);
14
15     [classes, prototypesHistory] = kmeans(data, K, true);
16     lastPrototypes = prototypesHistory(:, :, size(prototypesHistory, 3));
17     J20pp(i) = quantization_error(data, classes, lastPrototypes);
18 end
19
20 % Exercise b)
21 J20mean = mean(J20);
22 J20std = std(J20);
23
24 J20ppmean = mean(J20pp);
25 J20ppstd = std(J20pp);
26
27 % Exercise c)
28 [~, p] = ttest2(J20, J20pp, 'Vartype', 'unequal');
29
30 disp('Quantization error k-means:');
31 disp('Mean:');
32 disp(J20mean);
33 disp('Std:');
34 disp(J20std);
35
36 disp('Quantization error k-means++:');
37 disp('Mean:');
38 disp(J20ppmean);
39 disp('Std:');
40 disp(J20ppstd);
41
42 disp('p-value');
43 disp(p);

```

3.5 Implementation of the Batch Neural Gas algorithm

```

1 function [prototypes] = batchNG(Data, n, epochs, xdim, ydim)
2
3 % Batch Neural Gas
4 % Data contains data,
5 % n is the number of clusters,
6 % epoch the number of iterations,
7 % xdim and ydim are the dimensions to be plotted, default xdim=1,ydim=2
8
9 error(nargchk(3, 5, nargin)); % check the number of input arguments
10 if (nargin<4)
11     xdim=1; ydim=2; % default plot values
12 end;
13
14 [dlen, dim] = size(Data);
15
16 %prototypes = % small initial values
17 % % or
18 sbrace = @(x,y)(x{y});
19 fromfile = @(x)(sbrace(struct2cell(load(x)),1));
20 prototypes=fromfile('clusterCentroids.mat');
21
22 lambda0 = n/2; %initial neighborhood value
23 % lambda
24 lambda = lambda0 * (0.01/lambda0).^([0:(epochs-1)]/epochs);
25 % note: the lecture slides refer to this parameter as sigma^2
26 % instead of lambda

```

```

27
28 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
29 %% Action
30
31 for i=1:epochs,
32     D_prototypes = zeros(n,dim); % difference for vectors is initially zero
33     D_prototypes_av = zeros(n,1); % the same holds for the quotients
34
35     for j=1:dlen % consider all points at once for the batch update
36         % sample vector
37         x = Data(j,:); % sample vector
38         X = x(ones(n,1),:); % we'll need this
39
40         % neighborhood ranking
41
42         diffV = X - prototypes;
43         % We don't need the exact value of the distance, just a magnitude,
44         % so the distance squared is enough and more efficient to compute
45         distancesSq = diffV(:,1).^2 + diffV(:,2).^2;
46         [~, sortedPrototypeIndices] = sort(distancesSq);
47
48         % Ranking(i) is the index of the prototype in the sorted
49         % array
50         % Note that this is the number of prototypes with a distance
51         % smaller to the current data point because its sorted increasingly
52         ranking(sortedPrototypeIndices) = 1:n;
53
54
55         % DISTANCE!!!
56         % 1-BMU, 2-BMU, etc. (hint:sort)
57         %find ranking,h,H
58
59         % accumulate update
60         aux = exp(-ranking/lambda(i))';
61         D_prototypes = D_prototypes + repmat(aux, 1, 2).*X;
62         D_prototypes_av = D_prototypes_av + aux;
63     end
64
65     prototypes = D_prototypes./repmat(D_prototypes_av, 1, 2);
66
67     % track
68     if 1
69         fprintf(1, '%d / %d \r', i, epochs);
70         hold off
71         plot(Data(:,xdim), Data(:,ydim), 'bo', 'markersize', 3)
72         hold on
73         plot(prototypes(:,xdim), prototypes(:,ydim), 'r.', 'markersize', 10, 'linewidth', 3)
74         % write code to plot decision boundaries
75         voronoi(prototypes(:,1), prototypes(:,2));
76         %pause
77         %or
78         drawnow
79     end
80 end
81 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

```

3.6 Plotter to compare kmeans and batch neural gas

```

1 function kmeansbatchplotter()
2     function plotBatch(epochs)

```



```

3         batchNG(data, K, epochs);
4         xlabel('X');
5         ylabel('Y');
6         tt = sprintf('Batch Neural Gas clustering, %d epochs', epochs);
7         title(tt);
8     end
9
10    K = 100;
11
12    data = load('checkerboard');
13    data = data.checkerboard;
14
15    % Plot for k-means
16    figure(1);
17    [~, prototypesHistory] = kmeans(data, K, false, true);
18    iterations = size(prototypesHistory, 3);
19    lastPrototypes = prototypesHistory(:, :, iterations);
20    plot(data(:,1), data(:,2), 'bo', 'markersize', 3)
21    hold on
22    plot(lastPrototypes(:,1), lastPrototypes(:,2), 'r.', 'markersize', 10, 'linewidth', 3)
23    voronoi(lastPrototypes(:, 1), lastPrototypes(:, 2));
24    xlabel('X');
25    ylabel('Y');
26    tt = sprintf('%d-Means clustering', K);
27    title(tt);
28    hold off;
29    pause;
30
31    % Plot for batch neural 20
32    figure(2);
33    plotBatch(20);
34    pause;
35
36    % Plot for batch neural 100
37    figure(3);
38    plotBatch(100);
39    pause;
40
41    % Plot for batch neural 200
42    figure(4);
43    plotBatch(200);
44    pause;
45
46    % Plot for batch neural 500
47    figure(5);
48    plotBatch(500);
49
50    end

```