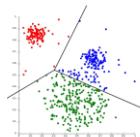


## 15-440

### Distributed Systems

### Kmeans

March 24, 2022  
Laila Elbeheiry



### K-Means at a High Level

- Clustering Algorithm

**Cluster analysis** or **clustering** is the task of grouping a set of objects in such a way that objects in the same group (called a **cluster**) are more similar (in some sense) to each other than to those in other groups (clusters).  
- Wikipedia

- Applications:

- Data mining
- Statistical data analysis: machine learning, pattern recognition, image analysis, information retrieval, and bioinformatics
- Visualization
- Detecting anomalies or outliers

What is "k"?

Carnegie Mellon University Qatar

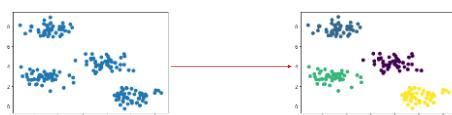


Image source:  
<https://towardsdatascience.com/k-means-clustering-explained-4528df86a120>

Carnegie Mellon University Qatar

### K-Means Explained

K-means is an iterative process that works by executing the following steps:

1. Select centroids (center of cluster) for each of the k clusters. The list of centroids can be selected by any method (e.g., randomly from the set of data points). It is usually better to pick centroids that are far apart.
2. Calculate the distance of all data points to the centroids.
3. Assign data points to the closest cluster.
4. Find the new centroids of each cluster by taking the mean of all data points in the cluster.
5. Repeat steps 2,3 and 4 until all points converge and cluster centers stop moving.

Let's see how it works!

Carnegie Mellon University Qatar

## K-Means in Project 3

You need to define a `distance` function and a `mean` function.

- How to calculate the distance between points in a 2D plane?
- How to calculate the distance between DNA strands?
- How to find the mean of points in a 2D plane?
- How to find the mean of DNA strands?

## Sequential Kmeans

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

Carnegie Mellon University Qatar

Initial centroids/means

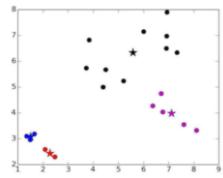
C0	P1
C1	P6
C2	P3
C3	P9
P1	
P2	
P3	
P4	
P5	
P6	
P7	
P8	
P9	
P10	
P11	
P12	
P13	
P14	
P15	
P16	

Carnegie Mellon University Qatar

Initial centroids/means

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

CO	P1
C1	P6
C2	P3
C3	P9



The blue and red starts are called unlucky centroids (\*). A poor choice of the initial centroids will take longer to converge or may result in bad clustering. You can handle this in:

1. Your data generators (generate first k points to be far apart and pick them in your implementation)
2. Try different sets of random centroids, and choose the best set.

Carnegie Mellon University Qatar

Initial centroids/means

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

CO	P1
C1	P6
C2	P3
C3	P9

Carnegie Mellon University Qatar

Initial centroids/means

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

4	CO	P1
9	C1	P6
2	C2	P3
10	C3	P9

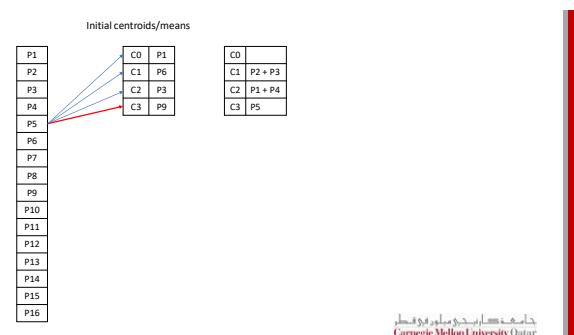
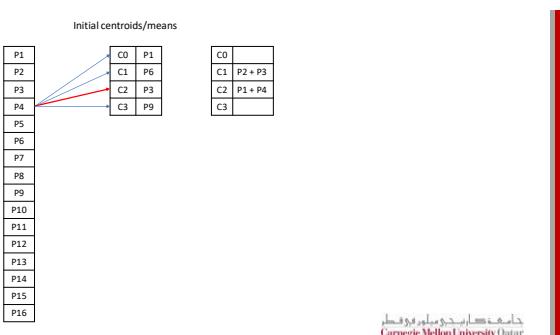
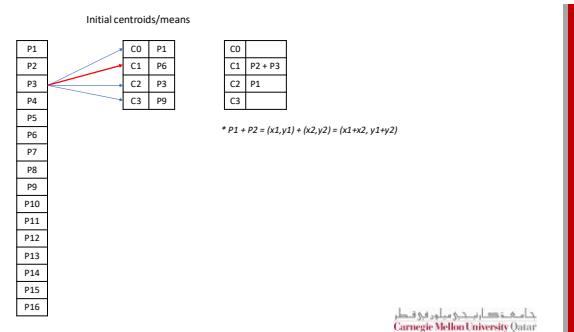
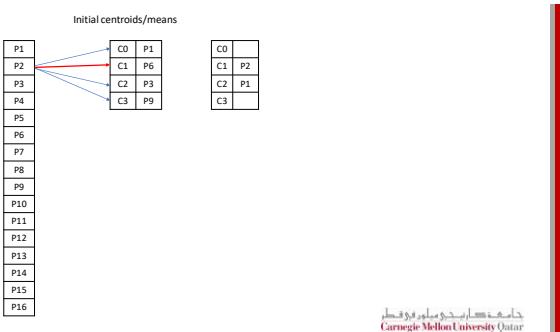
Carnegie Mellon University Qatar

Initial centroids/means

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

4	CO	P1
9	C1	P6
2	C2	P3
10	C3	P9

Carnegie Mellon University Qatar



Initial centroids/means

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

CO	P1
C1	P6
C2	P3
C3	P9

CO	P6 + P8 + P10 + P13 /4
C1	P2 + P3 + P7 + P11 /4
C2	P1 + P4 + P12 + P15 + P16 /5
C3	P5 + P9 + P14 /3

جامعة كارنيجي ميلون قطر | Carnegie Mellon University Qatar

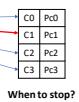
Centroids after iteration 1

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

CO	(P6 + P8 + P10 + P13)/4
C1	(P2 + P3 + P7 + P11)/4
C2	(P1 + P4 + P12 + P15 + P16)/5
C3	(P5 + P9 + P14)/3

\*  $\rho/N = (x/N, y/N)$

جامعة كارنيجي ميلون قطر | Carnegie Mellon University Qatar



When to stop?

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

جامعة كارنيجي ميلون قطر | Carnegie Mellon University Qatar

## Parallel K-Means

جامعة كارنيجي ميلون قطر | Carnegie Mellon University Qatar

How can we parallelize?

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

Carnegie Mellon University Qatar

How can we parallelize?

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

Carnegie Mellon University Qatar

How can we parallelize?

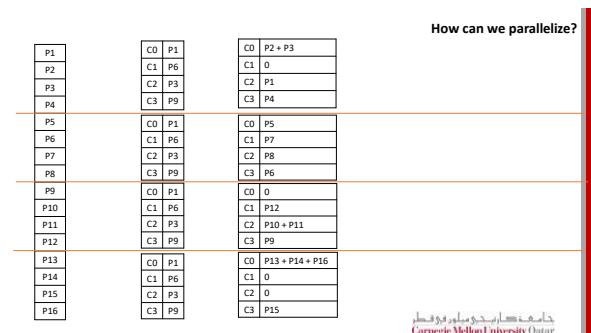
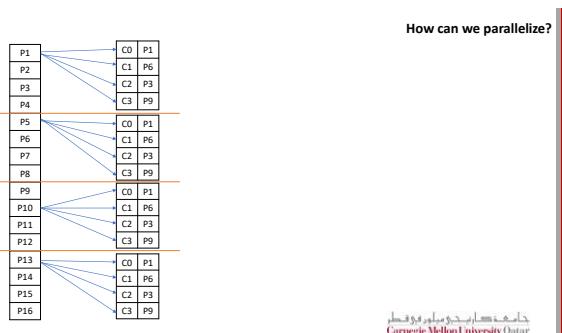
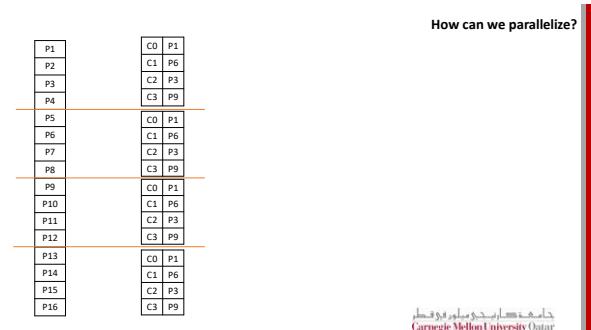
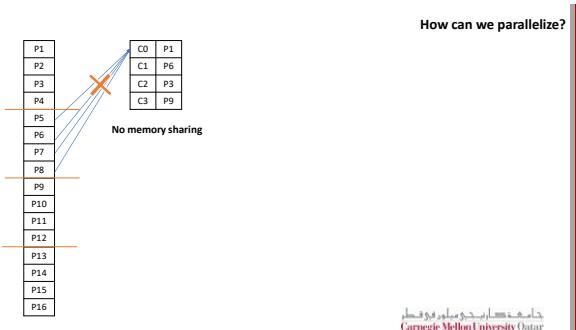
P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

Carnegie Mellon University Qatar

How can we parallelize?

P1
P2
P3
P4
P5
P6
P7
P8
P9
P10
P11
P12
P13
P14
P15
P16

Carnegie Mellon University Qatar



How can we parallelize?	
P1	C0   P1
P2	C1   P6
P3	C2   P3
P4	C3   P9
P5	C0   P1
P6	C1   P6
P7	C2   P3
P8	C3   P9
P9	C0   P1
P10	C1   P6
P11	C2   P3
P12	C3   P9
P13	C0   P1
P14	C1   P6
P15	C2   P3
P16	C3   P9
Carnegie Mellon University Qatar	

How can we parallelize?	
P1	C0   P1
P2	C1   P6
P3	C2   P3
P4	C3   P9
P5	C0   P5
P6	C1   P7
P7	C2   P8
P8	C3   P6
P9	C0   0
P10	C1   P12
P11	C2   P3
P12	C3   P9
P13	C0   P13 + P14 + P16
P14	C1   0
P15	C2   0
P16	C3   P15
Carnegie Mellon University Qatar	

How can we parallelize?	
P1	C0   P1
P2	C1   P6
P3	C2   P3
P4	C3   P9
P5	C0   P1
P6	C1   P6
P7	C2   P3
P8	C3   P9
P9	C0   P1
P10	C1   P6
P11	C2   P3
P12	C3   P9
P13	C0   P1
P14	C1   P6
P15	C2   P3
P16	C3   P9
Carnegie Mellon University Qatar	

How can we parallelize?	
P1	C0   P1
P2	C1   P6
P3	C2   P3
P4	C3   P9
P5	C0   P3 + P5 + P13 + P14 + P16 / 6
P6	C1   P7 + P12 / 2
P7	C2   P8 + P10 + P11 / 3
P8	C3   P6 + P9 + P15 / 3
P9	C0   P1
P10	C1   P6
P11	C2   P3
P12	C3   P9
P13	C0   P1
P14	C1   P6
P15	C2   P3
P16	C3   P9
Carnegie Mellon University Qatar	

How can we parallelize?

P1	
P2	
P3	
P4	
P5	
P6	
P7	
P8	
P9	
P10	
P11	
P12	
P13	
P14	
P15	
P16	

جامعة كارنيجي ميلون قطر

How can we parallelize?

P1	
P2	
P3	
P4	
P5	
P6	
P7	
P8	
P9	
P10	
P11	
P12	
P13	
P14	
P15	
P16	

جامعة كارنيجي ميلون قطر

## DNA stranding

ACTG
GTCA
SGGT
TAAA
ATAT

جامعة كارنيجي ميلون قطر

جامعة كارنيجي ميلون قطر

How to get the centroid of these DNA strands?

ACTG
GTCA
SGGT
TAAA
ATAT

How many repetitions of A in index 0 of all strands

A				
C				
G				
T				
Output strand				

Carnegie Mellon University Qatar

Carnegie Mellon University Qatar

ACTG
GTCA
SGGT
TAAA
ATAT

ACTG
GTCA
SGGT
TAAA
ATAT

A	2	1	2	1
C	0	1	1	0
G	1	1	1	1
T	1	2	1	2
Output strand				

Carnegie Mellon University Qatar

Carnegie Mellon University Qatar

ACTG
GTCA
SGGT
TAAA
ATAT

A	2	1	2	1
C	0	1	1	0
G	1	1	1	1
T	1	2	1	2
Output strand				

Get the mean or the median  
(sort the values and select the  
middle one)

جامعة كارنيجي ميلون قطر  
Carnegie Mellon University Qatar

ACTG
GTCA
SGGT
TAAA
ATAT

A	2	1	2	1
C	0	1	1	0
G	1	1	1	1
T	1	2	1	2
Output strand	T	G	C	A

جامعة كارنيجي ميلون قطر  
Carnegie Mellon University Qatar

ACTG
GTCA
SGGT
TAAA
ATAT

A	2	1	2	1
C	0	1	1	0
G	1	1	1	1
T	1	2	1	2
Output strand	T	G	C	A

جامعة كارنيجي ميلون قطر  
Carnegie Mellon University Qatar