

K-MEANS CLUSTERING

1. THE PROBLEM

13:17 - 13:29

2. ILLUSTRATION ON A 1-D EXAMPLE

16:15 - 16:42

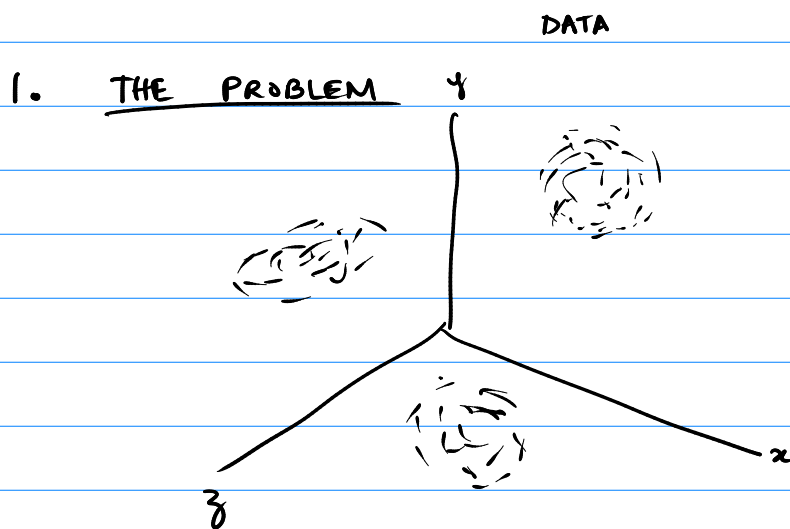
3. THE 1-D ALGORITHM

4. DISTORTION AND ITS MINIMIZATION BY K-MEANS

5. THE K-MEANS ALGORITHM FOR HIGHER-D DATA

6. "FAILURE" OF K-MEANS

7. EXAMPLES AND APPLICATIONS



FINDING CLUSTERS IN DATA

- GIVEN A , a matrix of DATA
- GIVEN k - NUMBER OF CLUSTERS
- FIND THE BEST CLUSTERS

→ APPLICATION: EG, MOVIE CLUSTERING USING RATINGS
- MANY OTHERS

→ PROBLEM VERY DIFFICULT (IN GENERAL): NP-HARD

→ A SOLUTION: K-MEANS CLUSTERING

→ OFTEN FINDS GOOD SOLUTIONS QUICKLY

→ NO GUARANTEE OF GOOD SOLUTION

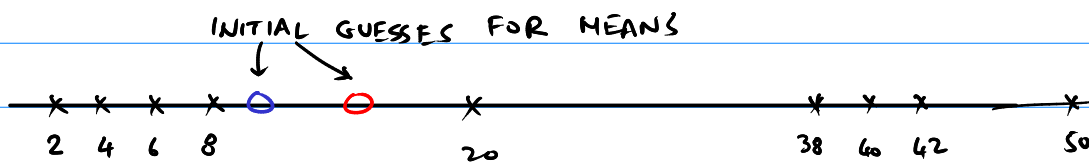
2. 1D-EXAMPLE OF K-MEANS



— 9 DATA POINTS

— GIVEN $k=2$ (find 2 best clusters)

— STEP 0: PICK $k=2$ GUESSES for CLUSTER CENTERS ("MEANS")



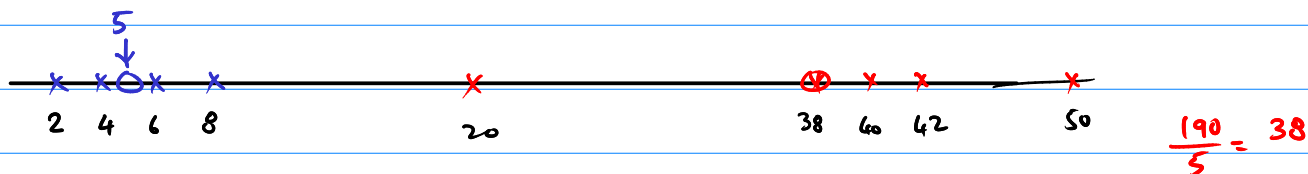
— ROUND 1

STEP 1: ASSIGN DATA TO NEAREST MEAN (CLUSTER)



STEP 2: UPDATE THE MEANS

— TO THE CENTROID/MEAN OF EACH CLUSTER'S DATA



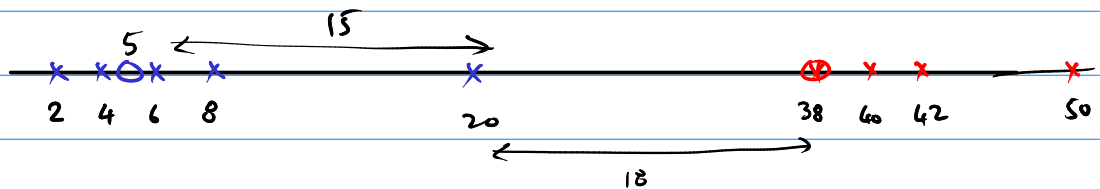
STEP 3 → HAVE CLUSTERS/MEANS CHANGED?

✓ → YES: GO BACK TO STEP 1, NEW ROUND

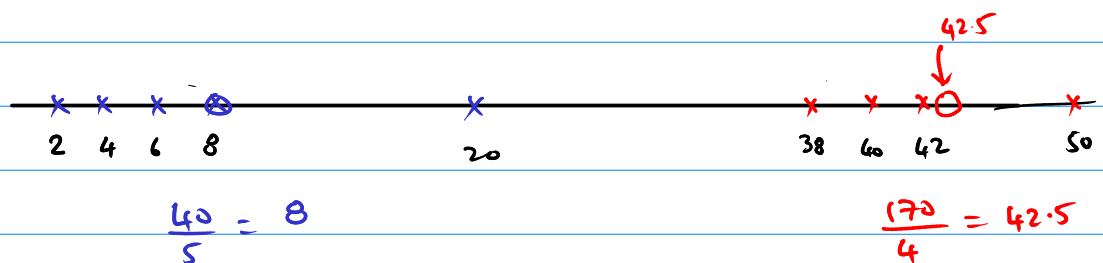
→ NO: DONE!

— ROUND 2

→ STEP 1: ASSIGN DATA TO NEAREST MEAN (CLUSTER)



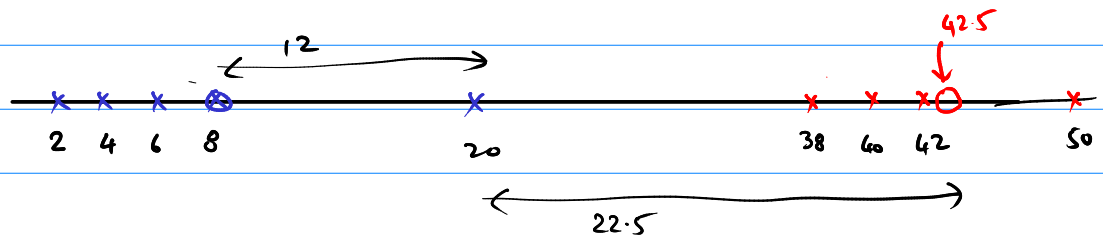
→ STEP 2: RE-CALCULATE CLUSTER MEANS



→ STEP 3: CLUSTERS CHANGED? YES. → ROUND 3, STEP 1

— ROUND 3:

→ STEP 1: RE-ASSIGN DATA TO NEAREST MEAN



→ NO CHANGE TO CLUSTERS

→ ⇒ (STEP 2) NO CHANGE TO MEANS.

→ DONE!

→ BLUE DATA IS ONE CLUSTER: MEAN = 8 ○

→ RED DATA IS 2nd CLUSTER: MEAN = 42.5 ○

3. THE ALGORITHM FOLLOWED ABOVE $S_1, \dots, S_k \leftarrow$ CLUSTER NAMES

A. PICK k GUESSES FOR MEANS: m_1, \dots, m_k

START ROUND:

B. ASSIGN DATA TO NEAREST MEAN (CLUSTER)

$$\rightarrow x_i \mapsto S_j \text{ s.t. } |x_i - m_j| \leq |x_i - m_\ell|, \ell = 1, \dots, k$$

\downarrow assigned to
 \downarrow CLOSEST MEAN FOR x_i

C. RE-CALCULATE MEANS

$$\rightarrow m_i = \frac{1}{|S_i|} \sum_{j=1}^{|S_i|} x_{ij}$$

\leftarrow j th data point in cluster i
 \downarrow # of elements in S_i (cardinality of S_i)

D. IF ANY CHANGES TO CLUSTERS/MEANS \rightarrow NEW ROUND, STEP B
 \rightarrow ELSE: DONE, k CLUSTERS ASSIGNED

4. THE DISTORTION METRIC AND ITS MINIMIZATION

- MEASURE OF HOW GOOD THE CLUSTER ASSIGNMENT IS
- THE SMALLER THE BETTER.

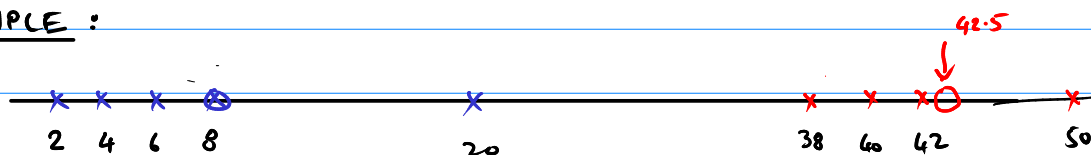
\rightarrow DISTORTION OF i th CLUSTER

$$D_i = \sum_{j=1}^{|S_i|} (x_{ij} - m_i)^2 \quad (\text{SUM OF SQUARES OF DISTANCES TO } m_i)$$

$$D = \sum_{i=1}^k D_i \quad (\text{SUM OF DISTORTIONS OF ALL CLUSTERS})$$

\rightarrow TOTAL DISTORTION

\rightarrow EXAMPLE:



$$D_1 = (2-8)^2 + (4-8)^2 + (6-8)^2 + (8-8)^2 + (20-8)^2 = 200$$

$$D_2 = (38-42.5)^2 + (40-42.5)^2 + (42-42.5)^2 + (50-42.5)^2 = 83$$

$$D = D_1 + D_2 = 283$$

MEASURE OF SPREAD OF CLUSTER = Variance $\times |S_i|$

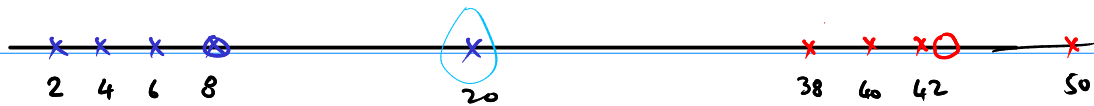
→ EACH STEP OF K-MEANS REDUCES DISTORTION

→ STEP 1: ASSIGN DATA TO NEAREST MEAN.

→ CHANGING THE ASSIGNMENT OF ANY ONE INCREASES D_i

↳ (i.e.) K-MEANS STEP 1 MINIMIZES D OVER CLUSTER ASSIGNMENTS
(WITH MEANS FIXED)

— EXAMPLE:



→ CONTRIBUTION OF 20 TO D_1 : $(20-8)^2 = 144$

→ LESS THAN $(20-42.5)^2 = 506.25 > 144$

→ MAKING 20 RED WILL INCREASE D_1

→ SAME FOR EVERY OTHER POINT

→ HENCE STEP 1 ASSIGNMENTS MINIMIZED D (FOR MEANS FIXED)

→ STEP 2: ASSIGN MEANS TO AVERAGE OF CLUSTER.

→ THIS CHOICE OF MEANS MINIMIZES D OVER ALL OTHER CHOICES (IF CLUSTERS ARE FIXED)

→ SUPPOSE p IS THE NEW MEAN INSTEAD.

$$D_i = \sum_{j=1}^{|S_i|} (x_j - p)^2$$

→ MINIMUM OF D_i WRT p : $\frac{\partial D_i}{\partial p} = 0$

$$\rightarrow \frac{\partial D_i}{\partial p} = \sum_{j=1}^{|S_i|} -2(x_j - p) = 0$$

$$\Rightarrow |S_i| p = \underbrace{\sum_{j=1}^{|S_i|} x_j}_{\text{this is } m_i}$$

→ HENCE $p = m_i$ minimizes D_i

→ TRUE FOR EVERY CLUSTER, HENCE $D = \sum_{i=1}^k D_i$ IS MINIMIZED.

5. THE GENERAL K-MEANS ALGORITHM (HIGHER-D DATA)

A. PICK k GUESSES FOR MEANS: $\vec{m}_1, \dots, \vec{m}_k \rightarrow q-d$

S_1 $S_k \leftarrow$ CLUSTER NAMES
↓ ↓
 \vec{m}_1 \vec{m}_k

→ START ROUND

B. ASSIGN DATA TO NEAREST MEAN (CLUSTER)

→ $\vec{x}_i \mapsto S_j$ s.t. $\|\vec{x}_i - \vec{m}_j\| \leq \|\vec{x}_i - \vec{m}_\ell\|, \ell = 1, \dots, k$

↓ assigned to ↓ CLOSEST MEAN FOR \vec{x}_i

C. RE-CALCULATE MEANS

→ $\vec{m}_i = \frac{1}{|S_i|} \sum_{j=1}^{|S_i|} \vec{x}_{ij}$

← j th data point in cluster i

↓ # of elements in S_i (cardinality of S_i)

D. IF ANY CHANGES TO CLUSTERS / MEANS → NEW ROUND, STEP B.

→ ELSE: STOP. k CLUSTERS ASSIGNED.

→ DISTORTION:

↳ DISTORTION OF i th CLUSTER

— $D_i =$ SUM OF SQUARES OF DISTANCES TO MEAN (FOR i th CLUSTER)

$$= \sum_{j=1}^{|S_i|} \|\vec{x}_{ij} - \vec{m}_i\|^2$$

— $D = \sum_{i=1}^k D_i$ (SUM OF DISTORTIONS OF ALL CLUSTERS)

↳ TOTAL DISTORTION

→ SAME MINIMIZATION PROPERTIES FOR STEPS 1 & 2.

— NUMERICAL DEMOS → SLIDES

6. "FAILURE" OF K-MEANS

→ K-MEANS DOESN'T ALWAYS DO A GOOD JOB



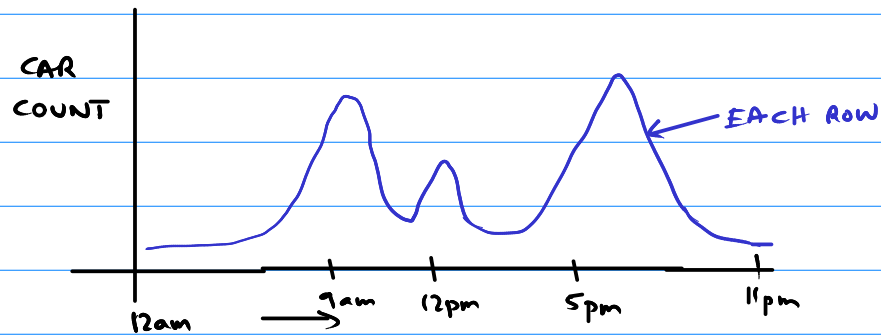
→ DEMOS OF BAD K-MEANS: → SLIDES

7. EXAMPLE APPLICATION: TRAFFIC PATTERNS

→ MONITOR TRAFFIC ACROSS AN INTERSECTION EVERY HOUR

	12-1am	1-2am	...	11 pm-12am
M	CAR COUNT	CAR COUNT		CAR COUNT
T			...	
W			...	
...			...	
F			...	
Sa			...	
Su			...	

— EACH ROW CAN BE PLOTTED:



→ TRY CLUSTERING THE DATA (A PAPER DID SO)

— 164 DAYS ; INTERSECTION IN S. CAROLINA

→ WITH $k=4$, THE CLUSTERS WERE:

M-Th, F, Sa, Su

→ (SHOW SLIDES)

→ HOW TO CHOOSE k ?

→ ELBOW CURVE METHOD:

— RUN K-MEANS FOR $k=2, 3, 4, \dots$

— PLOT DISTORTION AGAINST k

