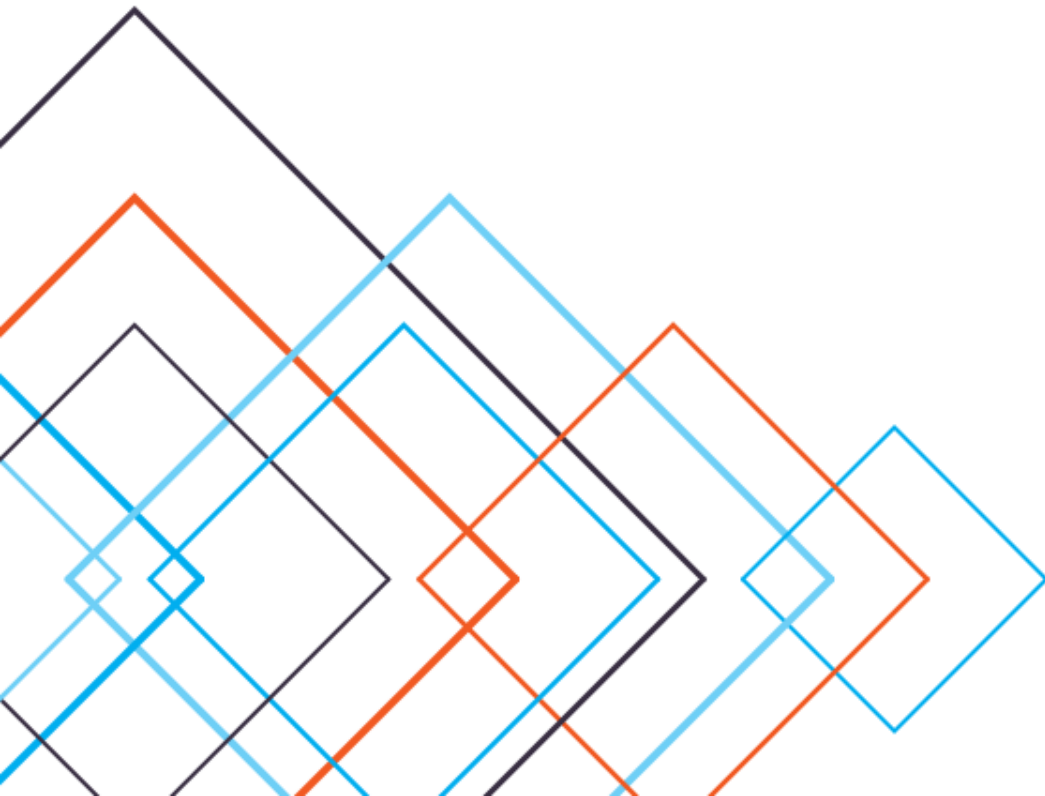# Distributed k-Means and k-Median Clustering on General Topologies



ISE Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
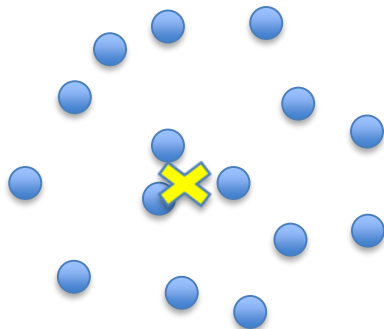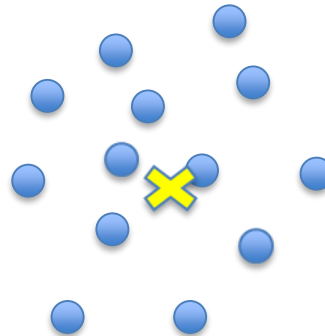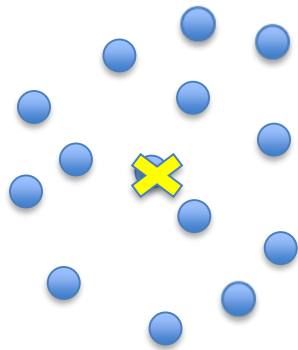
# Introduction

- Era of Big Data

- Data distributed over multiple locations
  - Distributed databases
  - Images and videos over networks
  - Sensors

- Centralized algorithms need to be scaled to distributed settings
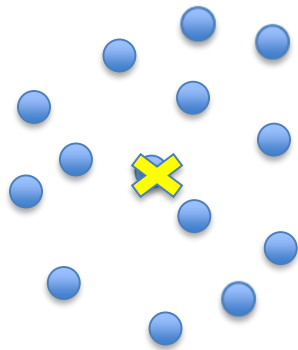
# Motivation

- k-Means and k-Medians are center based clustering algorithms

- Need a method to run k-means and k-medians on distributed data

- Communication cost should be low

- Run clustering algorithms on coreset

# What is a coreset/ ϵ-coreset?

- An ϵ-coreset is a weighted set of points whose cost on any set of centers is approximately the cost of the original data on those same centers up to accuracy ϵ.

Original data P with centers X

C = Cost( P , X )

Original data P with centers X

C = Cost( P , X )

P' = ϵ-coreset

ϵ-coreset P' with centers X

C = Cost( P , X )

P' = ϵ-coreset

C' = $\sum w(p')$Cost( p' , X )

ϵ-coreset P' with centers X

C = Cost( P , X )

P' = ∊-coreset

C' = $\sum w(p')$Cost( p' , X )

$(1 - \epsilon)$Cost( P , X ) $\leq \sum w(p')Cost(p', X) \leq (1 + \epsilon)$Cost( P , X )

∊-coreset P' with centers X

# What is a coreset/ ϵ-coreset?
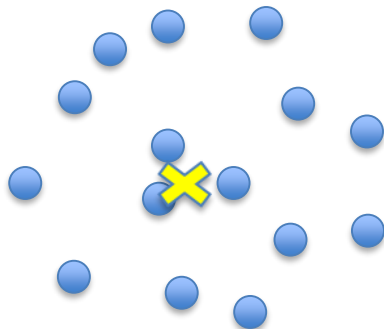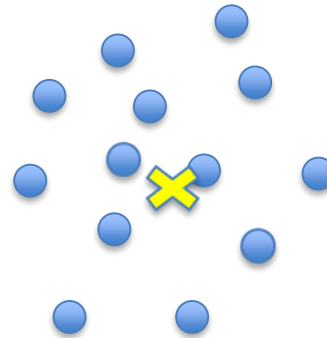
- An ϵ-coreset is a weighted set of points whose cost on any set of centers is approximately the cost of the original data on those same centers up to accuracy ϵ.
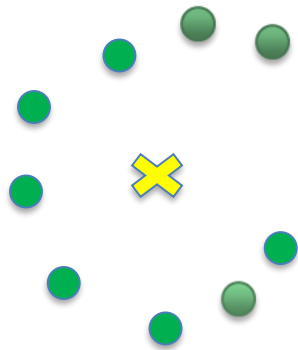
- Thus an approximate solution for the coreset is also an approximate solution for the original data.

# How are coresets constructed?

Original data P

**1) Compute constant approximate solution**

Original data P with centers X

**1) Compute constant approximate solution**

**2) Calculate cost of solution**

Original data P with centers X

Coreset P' with centers X

**1) Compute constant approximate solution**

**2) Calculate cost of solution**

**3) Sample points with probability proportional to their contributions to the cost.**

# **Methodology**

- The dataset is distributed on several nodes in an undirected connected graph.
- A local constant approximation solution is computed
- The total cost of local solution communicated to the other nodes in the graph using a message passing algorithm.
- Based on the cost of nodes and the contribution of the data towards their local solution, the local data are sampled to form local coresets

# Methodology

- The dataset is distributed on several nodes in an undirected connected graph.

- A local constant approximation solution is computed

- The total cost of local solution communicated to the other nodes in the graph using a message passing algorithm.

- Based on the cost of nodes and the contribution of the data towards their local solution, the local data are sampled to form local coresets

P4, X4

v6

P6, X6

v4

v3

P3, X3

v5

P5, X5

v2

P2, X2

v1

P1, X1

# Methodology

- The dataset is distributed on several nodes in an undirected connected graph.
- A local constant approximation solution is computed
- The total cost of local solution communicated to the other nodes in the graph using a message passing algorithm.
- Based on the cost of nodes and the contribution of the data towards their local solution, the local data are sampled to form local coresets

ISE Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Methodology

- The dataset is distributed on several nodes in an undirected connected graph.

- A local constant approximation solution is computed

- The total cost of local solution communicated to the other nodes in the graph using a message passing algorithm.

- Based on the cost of nodes and the contribution of the data towards their local solution, the local data are sampled to form local coresets

S4

v4

v6

S6

v3

S3

v5

S5

v2

S2

v1

S1

- The message passing algorithm is called once again to combine the local coresets.

- If each node constructs an ∈-coreset on its local data, then the union of these local coresets is an ∈-coreset for the entire data.

- The constant approximation algorithm is applied to this global coreset to obtain the final result.

- The message passing algorithm is called once again to combine the local coresets.

- If each node constructs an $\epsilon$-coreset on its local data, then the union of these local coresets is an $\epsilon$-coreset for the entire data.

- The constant approximation algorithm is applied to this global coreset to obtain the final result.

- The message passing algorithm is called once again to combine the local coresets.

- If each node constructs an ϵ-coreset on its local data, then the union of these local coresets is an ϵ-coreset for the entire data.

- The constant approximation algorithm is applied to this global coreset to obtain the final result.

# Comparison to Other Coreset Algorithms:

- Zhang et. al.[1] is limited to rooted trees.
  - The communication cost to develop coresets is very high.
  - Although it is possible to find a spanning tree, when the graph has large diameter every tree has large height which greatly increases the size of coresets which must be communicated across the lower levels of the tree.
- D. Feldman et. al.[2] ignores the communication cost and merges coresets in a parallel computation model.

# Comparison to Other Coreset Algorithms:

- Zhang et. al.[1] is limited to rooted trees.
  - The communication cost to develop coresets is very high.
  - Although it is possible to find a spanning tree, when the graph has large diameter every tree has large height which greatly increases the size of coresets which must be communicated across the lower levels of the tree.
- D. Feldman et. al.[2] ignores the communication cost and merges coresets in a parallel computation model.

# Technical Background

# **Preliminaries**

- Consider two points $p, q \in R^d$; $d(p, q)$ is the Euclidean distance between p and q

- k-means goal: Find set of k centers $X = \{x_1, x_2, \ldots, x_k\}$ which minimize the k-means cost of data $P \subseteq R^d$

- k-means cost = $Cost(P, X) = \sum_{p \in P} d(p, X)^2$
  where $d(p, X) = \min_{x \in X} d(p, X)$

- If P is a weighted dataset with a weight function w, then k-means cost = $Cost(P, X) = \sum_{p \in P} w(p) d(p, X)^2$

- For distributed clustering, we consider Set of n nodes $V = \{v_i, 1 \leq i \leq n\}$ which communicate on undirected connected graph $G = (V, E)$, with $m = |E|$ edges.

- Each node $v_i$ contains local set of data points $P_i$ and the global dataset is $\mathrm{P} = \bigcup_{i=1}^{n} P_i$

- $(v_i, v_j) \in E$ indicates that $v_i \ and \ v_j$ can communicate with each other.

- Communication cost = number of points transmitted. For simplicity, assume that there is no latency in communication.

# Distributed coreset construction

**Input:** Data (distributed across different nodes)



- $P_i$ : Set of points in node i
- Round 1 is performed on each node followed by Round 2.
- Only cost values are communicated.

**Round 1**

Compute centers for local data

Compute constant approximation

Communicate costs to nodes

**Round 2**

Compute number of points to be sampled

Set weights on local data

Random Sampling of points to coreset

Inclusion of local centers to coreset

Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Round 1

**Step 1:** Computing centers for local data



- Perform k-center clustering on node data (Lloyd's Algorithm)
- Output: Set of centers $B_i$ for $P_i$

| Round 1 |
| --- |
| Compute centers for local data |
| Compute constant approximation |
| Communicate costs to nodes |

| Round 2 |
| --- |
| Compute number of points to be sampled |
| Set weights on local data |
| Random Sampling of points to coreset |
| Inclusion of local centers to coreset |

# Round 1

**Step 2:** Computing constant approximation



$$Cost(P_i, B_i) = \begin{cases} \sum_{p \in P_i} d(p, B_i) \ for \ k-medians \\ \sum_{p \in P_i} d(p, B_i)^2 \ for \ k-means \end{cases}$$

where $d(p, B_i) = \min_{b \in B_i} d(p, b)$

| Round 1 |
|---|
| Compute centers for local data |
| Compute constant approximation |
| Communicate costs to nodes |

| Round 2 |
|---|
| Compute number of points to be sampled |
| Set weights on local data |
| Random Sampling of points to coreset |
| Inclusion of local centers to coreset |

**ISE** Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Round 1

**Step 3:** Communicating cost to nodes



- $Cost(P_i, B_i)$ is communicated to each node.
- *Message_Passing()* used for cost communication

| Round 1 |
| --- |
| Compute centers for local data |
| Compute constant approximation |
| Communicate costs to nodes |

| Round 2 |
| --- |
| Compute number of points to be sampled |
| Set weights on local data |
| Random Sampling of points to coreset |
| Inclusion of local centers to coreset |

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Round 2

**Step 1:** Computing # of points to be sampled



- $t$ = Global coreset size
- Number of points to be sampled on node i, $t_i = \dfrac{t \, cost(P_i, B_i)}{\sum_{j=1}^{n} cost(P_j, B_j)}$
- Proportional to total cost contribution

**Round 1**

Compute centers for local data

Compute constant approximation

Communicate costs to nodes

**Round 2**

Compute number of points to be sampled

Set weights on local data

Random Sampling of points to coreset

Inclusion of local centers to coreset

Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Round 2

**Step 2:** Setting of weights on local data



- Weight of data point p, $m_p = 2cost(p, B_i) \, \forall \, p \in P_i$
- Proportional to distance of point from corresponding local center

**Round 1**

Compute centers for local data

Compute constant approximation

Communicate costs to nodes

**Round 2**

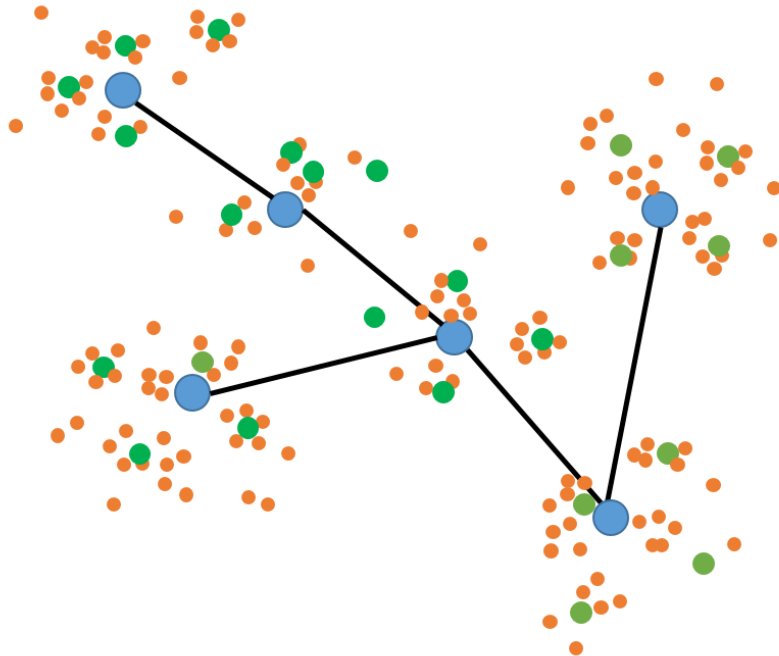Compute number of points to be sampled

Set weights on local data

Random Sampling of points to coreset

Inclusion of local centers to coreset

Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Round 2

**Step 3:** Random sampling of points to coreset



- Non-uniform random sample $S_i$ of $t_i$ points from $P_i$, where for every $q \in Pi$, we have

$$\Pr[q = p] = \frac{m_p}{\sum_{z \in P_i} m_z}$$

- Weight on q, $w_q = \dfrac{\sum_i \sum_{z \in P_i} m_z}{m_q \, t}$ each $q \in Si$



Round 1

Compute centers for local data

Compute constant approximation

Communicate costs to nodes
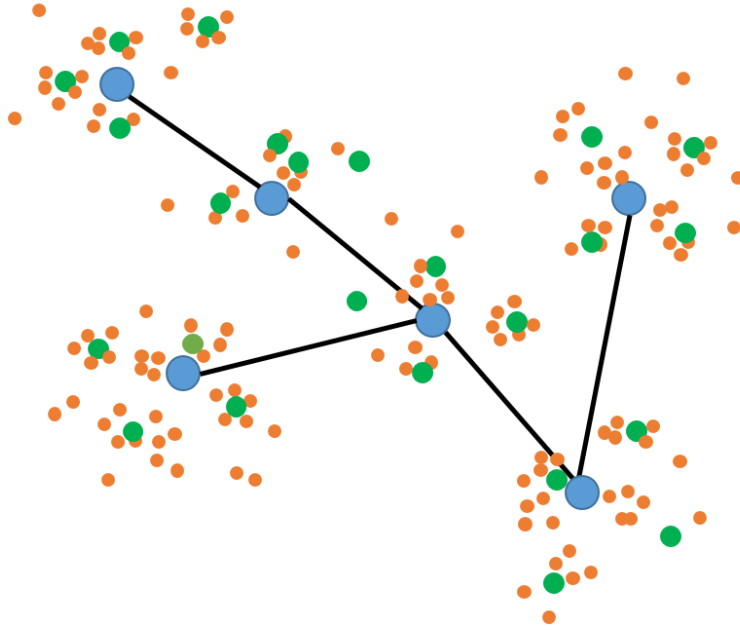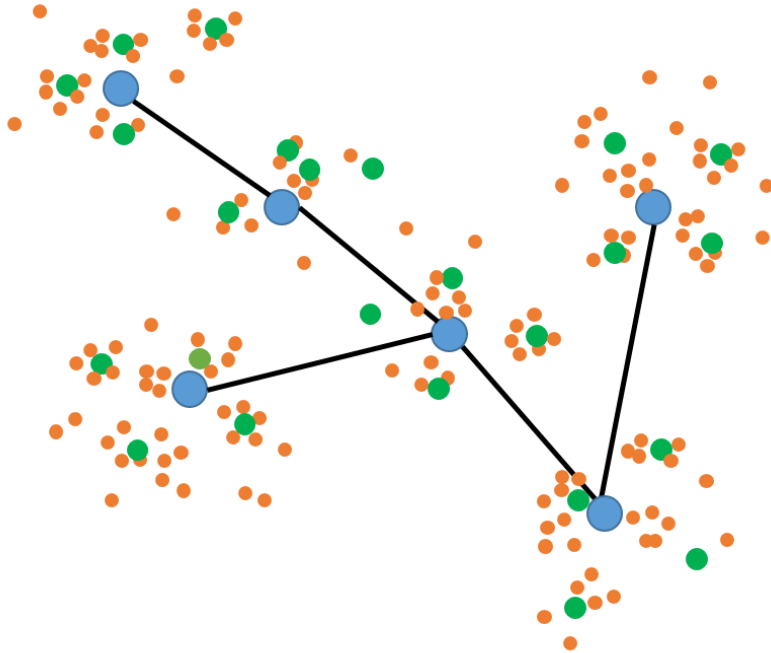
Round 2

Compute number of points to be sampled

Set weights on local data

Random Sampling of points to coreset

Inclusion of local centers to coreset

# Round 2

**Step 4:** Including local centers to coreset



- Local k-centers ($B_i$) are included to the coreset
- Weight on b, $w_b = |P_b| - \sum_{k \in P_b \cap S} w_k$, where
  $P_b = \{ p \in P_i : d(p, b) = d(p, B_i) \}$



**Round 1**

Compute centers for local data

Compute constant approximation

Communicate costs to nodes

**Round 2**

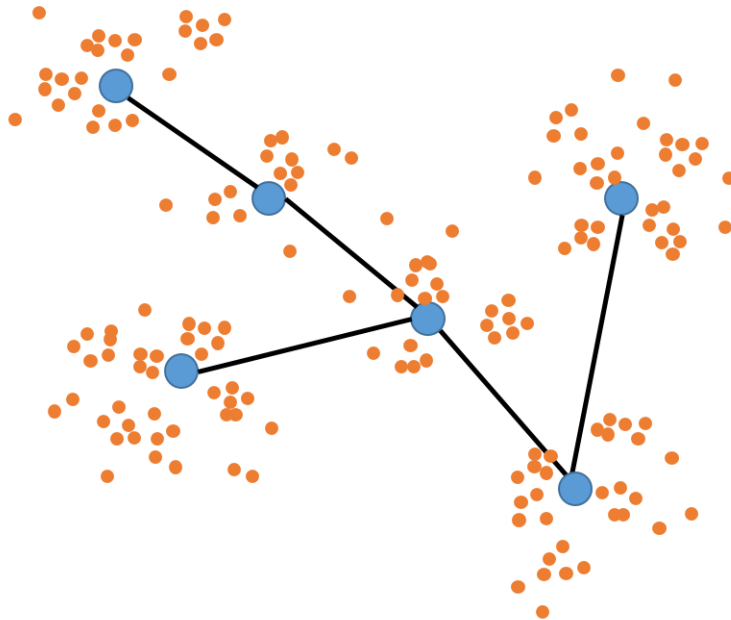Compute number of points to be sampled

Set weights on local data

Random Sampling of points to coreset

Inclusion of local centers to coreset

# Distributed coreset construction

**Input:** Data (distributed across different nodes)



**Output:** Weighted Coreset



**Round 1**

Compute centers for local data

Compute constant approximation

Communicate costs to nodes

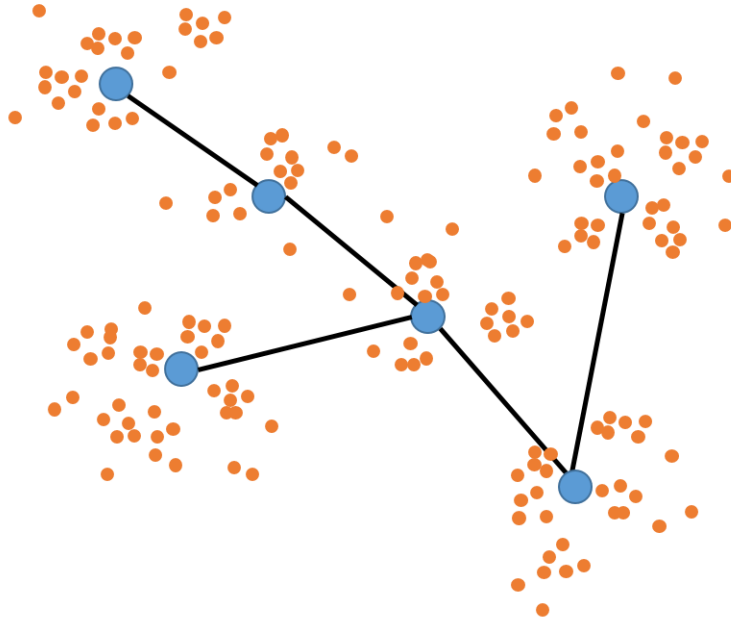**Round 2**

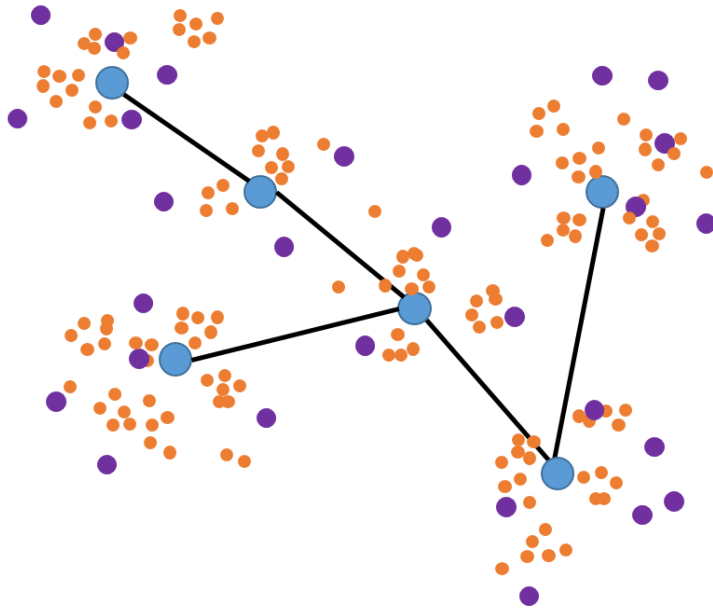Compute number of points to be sampled

Set weights on local data

Random Sampling of points to coreset

Inclusion of local centers to coreset

## Objective:

- Reduce the total communication cost
- Preserve the theoretical guarantees for approximate clustering cost

## Theorem 1:

There exists an algorithm that with probability at least $1 - \delta$, the output is an $\epsilon$-coreset for the global dataset. The communication cost is $O(mn)$ and the number of points in the coresets (*coreset size*) are:

- **k- means:** $O\left(\frac{1}{\epsilon^4}\left(kd + log\left(\frac{1}{\delta}\right)\right) + nk\ log\left(\frac{nk}{\delta}\right)\right)$

- **k-medians:** $O\left(\frac{1}{\epsilon^2}\left(kd + log\left(\frac{1}{\delta}\right)\right) + nk\right)$

# Key Parameters

- **Sampling Probability:** Probability with which data points are chosen to be candidates for coresets

$$\Pr[q = p] = \frac{m_p}{\sum_{z \in P} m_z}$$

- **Weights on coreset points:** The weights on coreset points are:

$$w_q = \frac{\sum_i \sum_{z \in P_i} m_z}{m_q \, |S|} \text{ for q} \in S \text{ and } w_b = |P_b| - \sum_{k \in P_b \cap S} w_k \text{ for } b \in B$$

- **Function applied on the points**: A suitably defined function $f(p)$ is selected. The function nature varies for k-means and k-medians.

# Validity of Sampling probability and Weights

**Claim:** For a given function $f(p)$, the sampled coreset points are valid approximations of the original data

**To Prove:** $E\left[\sum_{q \in S} w_q f(q)\right] = \sum_{p \in P} f(p)$ for chosen probabilities $\Pr[q = p]$ and weights $w_p$

**Proof:**

$$E\left[\sum_{q \in S} w_q f(q)\right] = \sum_{q \in S} E[w_q f(q)] = \sum_{q \in S} \sum_{p \in P} \Pr[q = p] w_p f(p)$$

$$= \sum_{q \in S} \sum_{p \in P} \frac{m_p}{\sum_{z \in P} m_z} \frac{\sum_{z \in P} m_z}{m_p |S|} f(p)$$

$$= \sum_{q \in S} \sum_{p \in P} \frac{1}{|S|} f(p) = \sum_{p \in P} f(p)$$

# Validity of *f(p)*

## Lemma 1:

Fix a set F of functions $f\colon P \to \mathbb{R}_+$. Let S be an i.i.d. sample drawn from P according to $\{m_p \colon p \in P\}$. For $p \in P$ and $q \in S$, we have $q = p$ with probability $\frac{m_p}{\sum_{z \in P} m_z}$. Also, define $w_p = \frac{\sum_{z \in P} m_z}{m_p |S|}$. If for a sufficiently large c, we have $|S| \geq \frac{c}{\epsilon^2} \left( \dim(F, P) + \log \frac{1}{\delta} \right)$, then with a probability at least 1 - δ ,

$$\forall\, f \in F \colon \left| \sum_{p \in P} f(p) - \sum_{q \in S} w_q f(q) \right| \leq \epsilon \left( \sum_{p \in P} m_p \right) \left( max_{p \in P} \frac{f(p)}{m_p} \right)$$

- Set $m_p = max_{f \in F} f(p)$. Bounds RHS to $\epsilon \left( \sum_{p \in P} m_p \right)$
- Natural Approach: Set $f_x(p) \coloneqq cost(p, x)$
- Complication: Suitable upper bound not available for $cost(p, x)$
- Alternative Approach: Set $f_x(p) \coloneqq cost(p, x) - cost(b_p, x) + cost(p, b_p)$

p — x

b_p

$$0 \leq f_x(p) \leq 2cost(p, b_p)$$

# Proof of Theorem 1 for k-medians

**Outline:** Given that the sampling probabilities, weights on points and $f(p)$ as defined previously, algorithm outputs an $\epsilon$-coreset.

- Set $f_x(p) := d(p,x) - d(b_p, x) + d(p, b_p)$. Then, $m_p = 2d(p, b_p)$ and the conditions stated in Lemma 1 are satisfied.

- Then,

$$D = \left| \sum_{p \in P} f_x(p) - \sum_{q \in S} w_q f_x(q) \right| \leq O(\epsilon) \sum_{p \in P} d(p, x)$$

- Substituting the definition of $f_x(\circ)$ in the expression for D,

$$D = \left| \sum_{p \in P} [d(p,x) - d(b_p, x) + m_p] - \sum_{q \in S} w_q [d(p,x) - d(b_p, x) + m_q] \right|$$

- On further simplification

$$D = \left| \sum_{p \in P} d(p, x) - \sum_{q \in S \cup B} w_q d(q, x) \right| \leq O(\epsilon) \sum_{p \in P} d(p, x)$$

which is guarantee that the output of algorithm 1 is an $\epsilon$-coreset

# Proof of Theorem 1 for k-means

- **Complication:**

  Difference term $\left|cost(p,x) - cost(b_p,x)\right| = \left|d(p,x)^2 - d(b_p,x)^2\right|$ is not bounded.
  Conditions of Lemma 1 are not met.

- **Fix:** Two categories of points:
- ❑ **Good points** (costs approximately satisfy the triangle inequality)
- ❑ **Bad points**

## GOOD POINTS:

- – Define the good points with centers x as
- – $G(x) = \{p \in P : \left|cost(p,x) - cost(b_p,x)\right| \leq \Delta_p\}$

  where $\Delta_p = \dfrac{cost(b_p,x)}{\epsilon}$
- – These difference between the costs can be bound as before.

**BAD PONTS:**

For bad points, it is proven that the difference in cost may be larger than $\frac{cost(b_p,x)}{\epsilon}$ but will be $O(\epsilon \min\{cost(p,x), cost(b_p,x)\})$.

- Hence $f_x(p)$ is defined only over good points as follows:

$$f_x(p) = \begin{cases} cost(p,x) - cost(b_p,x) + \Delta_p & if \ p \in G(x) \\ 0 & otherwise \end{cases}$$

- Then, $D = \left| \sum_{p \in P} cost(p,x) - \sum_{q \in S \cup B} w_q cost(q,x) \right|$ can be decomposed into three terms

$$\sum_{p \in P} f_x(p) - \sum_{q \in S} w_q f_x(q) \qquad (1)$$
$$+ \sum_{p \in P \backslash G(x)} [cost(p,x) - cost(b_p,x) + \Delta_p] \qquad (2)$$
$$- \sum_{p \in P \backslash G(x)} w_q [cost(q,x) - cost(b_q,x) + \Delta_q] \qquad (3)$$

Industrial and Enterprise Systems Engineering

# Proof of Theorem 1 for k-means (cont.)

- It can be showed that each of these terms is bounded by $O(\epsilon)cost(P, x)$

- Sum of three terms is also bounded by $O(\epsilon)cost(P, x)$

- Then,

$$D = \left| \sum_{p \in P} cost(p, x) - \sum_{q \in S \cup B} w_q cost(q, x) \right| \leq O(\epsilon)cost(P, x)$$

- Hence, by choosing an appropriate $\epsilon$, the result is established.

# Pseudo code

**Git Link**: https://github.com/harshadrai/Distributed_kMeans

# Distributed coreset construction

**Input:** Data (distributed across different nodes)



**Output:** Weighted Coreset



| Round 1 |
|---|
| Compute centers for local data |
| Compute constant approximation |
| Communicate costs to nodes |

| Round 2 |
|---|
| Compute number of points to be sampled |
| Set weights on local data |
| Random Sampling of points to coreset |
| Inclusion of local centers to coreset |

ISE
Industrial and Enterprise

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Message Passing

**Input:**

**Message** to communicate
**Neighbors** to whom we wish to communicate the message



$I_i$ - Message at node $i$
$N_i$ - Neighbors of node $i$

$\implies N_i$ - Node 2

---

INPUT: Message, Neighbors

**Message Passing($I_i$, $N_i$)**

Define $R_i=\{I_i\}$ as information recieved

Send $I_i$ to all neighbors $N_i$

While $R_i \neq \{I_j : 1 \leq j \leq n\}$

NO

If $I_j \notin R_i$

YES

$R_i = R_i \cup \{I_j\}$

Send $I_j$ to all neighbors $N_j$

**Step 1:**

Define $R_i = \{I_i\}$ as the information received



Node $i$ with message $I_i$

$I_i$ - Message at node $i$
$N_i$ - Neighbors of node $i$

$\implies N_i$ - Node 2
$\qquad R_i = \{I_i\}$

---

INPUT: Message, Neighbors

**Message Passing($I_i$, $N_i$)**

Define $R_i = \{I_i\}$ as information recieved

Send $I_i$ to all neighbors $N_i$

While $R_i \neq \{I_j : 1 \leq j \leq n\}$

NO

If $I_j \notin R_i$

YES

$R_i = R_i \cup \{I_j\}$

Send $I_j$ to all neighbors $N_j$

---

**Step 2:**

Send $I_i$ to all its' neighbors $N_i$



$I_i$
1
Node $i$ with message $I_i$
2 Message $I_i$ received
3

INPUT: Message, Neighbors

Message Passing($I_i$, $N_i$)

Define $R_i$={$I_i$} as information recieved

Send $I_i$ to all neighbors $N_i$

While $R_i \neq \{I_j : 1 \leq j \leq n\}$

NO
If $I_j \notin R_i$
YES

$R_i = R_i \cup \{I_j\}$

Send $I_j$ to all neighbors $N_j$

ISE Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

**Step 3:**

Check if information received $R_i$ is same as information present at all the nodes
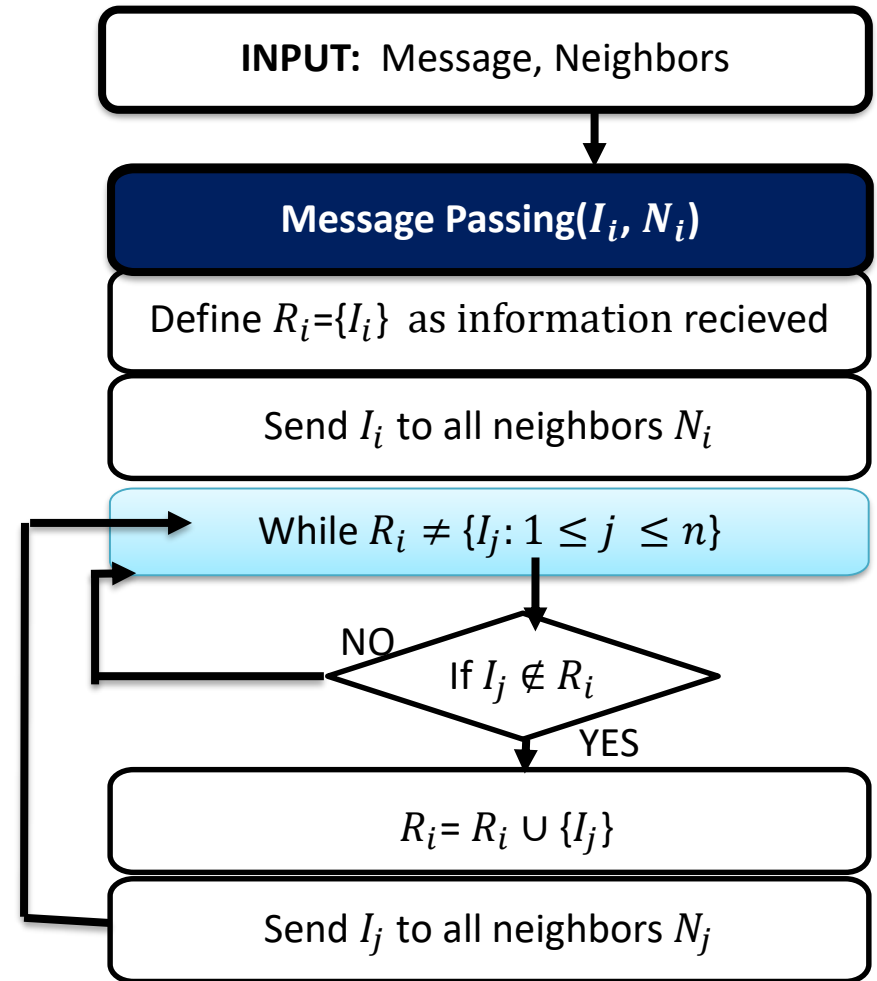
**YES:** STOP

**NO:** Continue to step 4

INPUT: Message, Neighbors

Message Passing($I_i$, $N_i$)

Define $R_i=\{I_i\}$ as information recieved

Send $I_i$ to all neighbors $N_i$

While $R_i \neq \{I_j : 1 \leq j \leq n\}$

If $I_j \notin R_i$

NO

YES

$R_i = R_i \cup \{I_j\}$
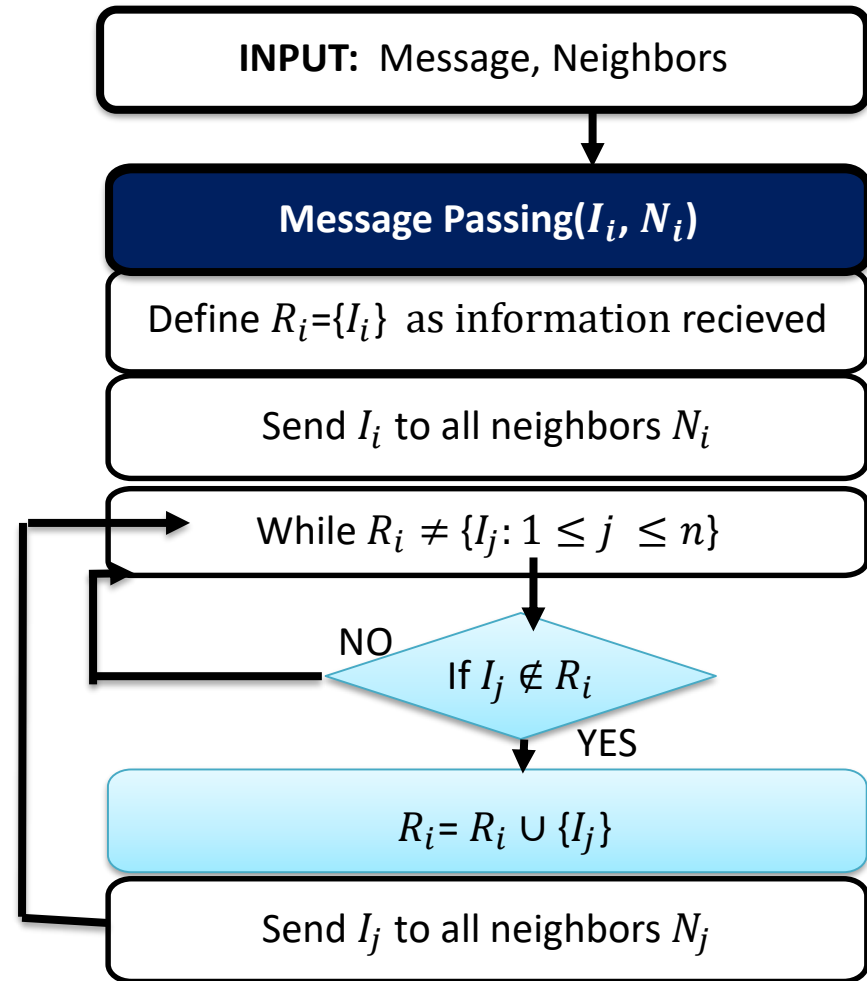
Send $I_j$ to all neighbors $N_j$

**Step 4:**

Check if information $I_j \notin R_i$

**YES:** Extend set of information received by $I_j$

$\quad R_i = R_i \cup \{I_j\}$

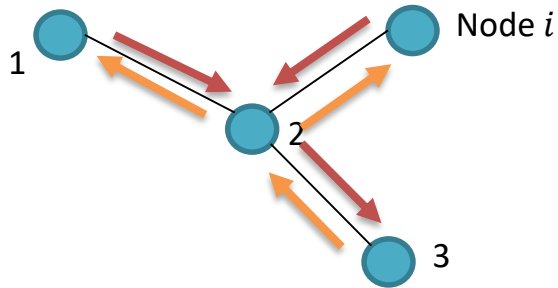**NO:** Continue to step 3



INPUT: Message, Neighbors

Message Passing($I_i$, $N_i$)

Define $R_i = \{I_i\}$ as information recieved

Send $I_i$ to all neighbors $N_i$

While $R_i \neq \{I_j : 1 \leq j \leq n\}$

If $I_j \notin R_i$

NO

YES

$R_i = R_i \cup \{I_j\}$

Send $I_j$ to all neighbors $N_j$

**Step 5:**

Send $I_j$ to all its' neighbors $N_j$



INPUT: Message, Neighbors

Message Passing($I_i$, $N_i$)

Define $R_i$={$I_i$} as information recieved

Send $I_i$ to all neighbors $N_i$

While $R_i \neq \{I_j : 1 \leq j \leq n\}$

NO

If $I_j \notin R_i$

YES

$R_i = R_i \cup \{I_j\}$

Send $I_j$ to all neighbors $N_j$

ISE Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

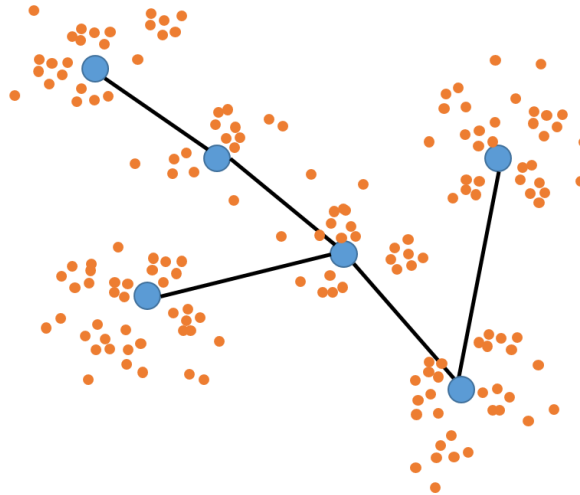**Output:**

Communication done among all the nodes


Node $i$

- As a result, information is shared globally.

- Communication cost at each step = Size of the message

# Distributed clustering on graph

**Input:**

**Data** (distributed among nodes)

$A_\alpha$ - $\alpha$ approximation clustering algorithm (**Lloyd's Algorithm**)



- $P_i$ : Set of points in node $i$
- $N_i$ : Neighbors at node $i$

**Round 1**

For each node $i$

Get $\varepsilon$ local coreset $D_i$ using Distributed Coreset construction algorithm

**Round 2**

For each node $i$

Call Message_Parsing ( $D_i$, $N_i$ )
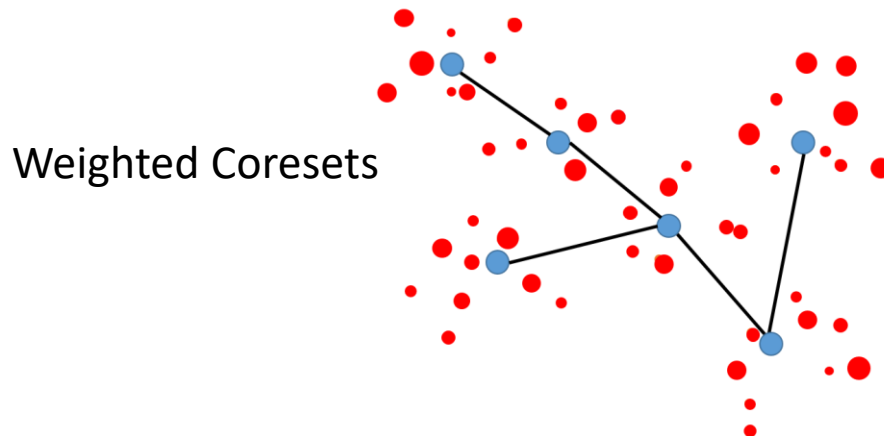
$X = A_\alpha$ ($\cup D_j$)

# Round 1

**Step 1:**
Apply Distributed Coreset construction algorithm on data

**Output:**
Local weighted coresets $D_i$ on each node $i$

Weighted Coresets



Round 1

For each node $i$

Get $\varepsilon$ local coreset $D_i$ using Distributed Coreset construction algorithm

Round 2
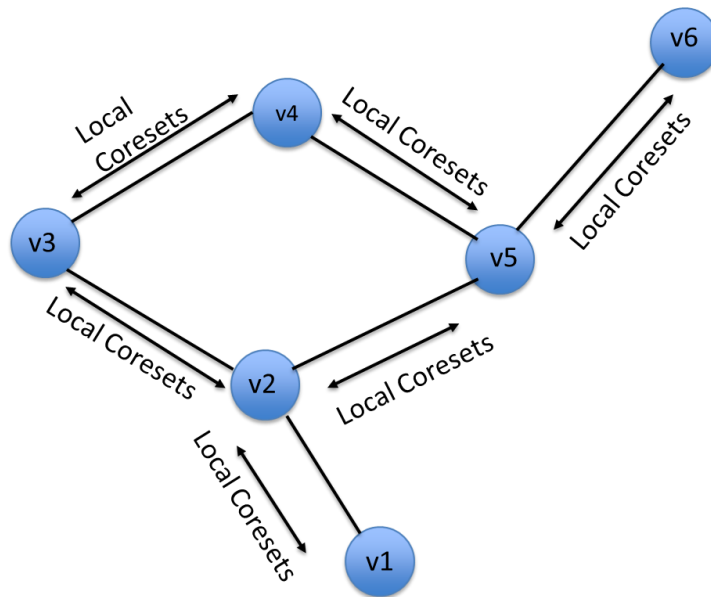
For each node $i$

Call Message_Parsing ( $D_i$, $N_i$ )

X = $A_\alpha$ ($\cup$ $D_j$)

# Round 2

**Step 1:**

Distributed coresets are communicated to each node
Using Message Passing algorithm



Round 1

For each node $i$

Get $\varepsilon$ local coreset $D_i$ using Distributed Coreset construction algorithm

Round 2

For each node $i$

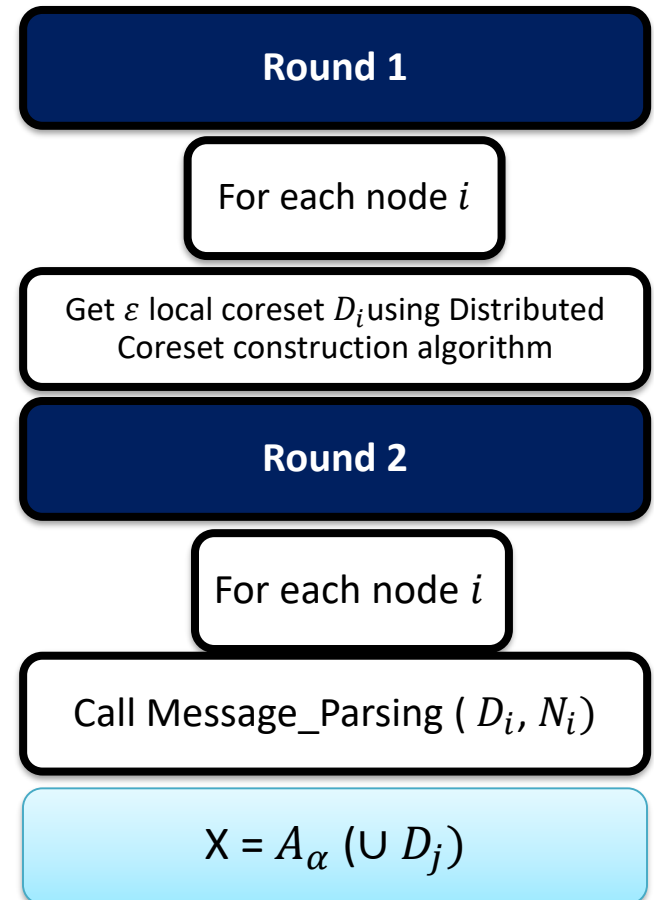Call Message_Parsing ( $D_i$, $N_i$ )

X = $A_\alpha$ ($\cup D_j$)

# Round 2

**Step 2:**

Lloyd's clustering algorithm is applied to compute the centers for the global dataset

**Output:**

Centers of the Global dataset

---

**Round 1**

For each node $i$

Get $\varepsilon$ local coreset $D_i$ using Distributed Coreset construction algorithm

**Round 2**

For each node $i$

Call Message_Parsing ( $D_i$, $N_i$ )

$X = A_\alpha\ (\cup\ D_j)$

---

ISE  Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Random Graph

Each edge is generated independently with probability p.

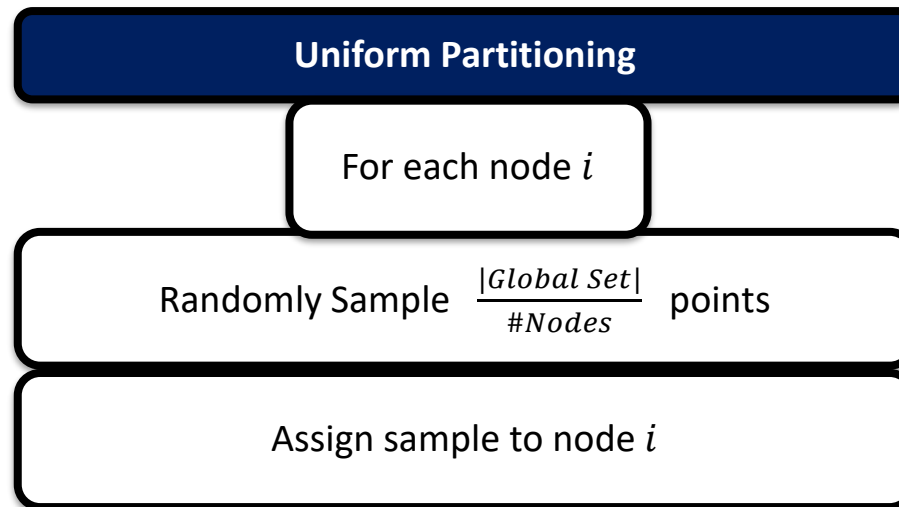For experiments, we used p = 0.3 and 10 nodes.

# Preferential Graph

Preferential graphs are generated using the preferential attachment process.

Additional edges are added continuously to the graph that are distributed among the nodes as an increasing function of the number of edges the nodes already have.
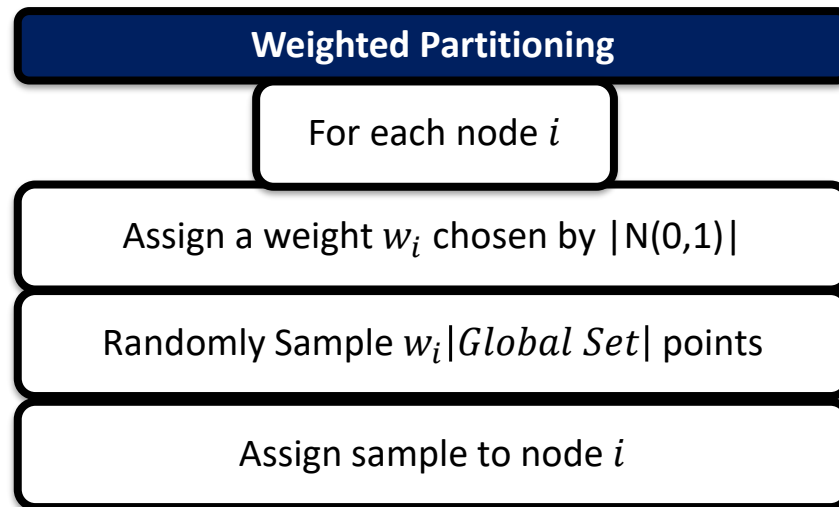
# Uniform Partitioning

Each data point from global dataset is assigned to a local site with equal probability.

<div style="border:1px solid #000;">

**Uniform Partitioning**

For each node $i$

Randomly Sample $\frac{|Global\ Set|}{\#Nodes}$ points

Assign sample to node $i$

</div>

**P(assigning a global data point to a node $i$)** $= \dfrac{1}{No.\ of\ nodes}$

# Weighted Partitioning

Each data point from global dataset is assigned to a local site with probability proportional to the nodes' weight.

| Weighted Partitioning |
|---|
| For each node $i$ |
| Assign a weight $w_i$ chosen by $|\text{N(0,1)}|$ |
| Randomly Sample $w_i|Global\ Set|$ points |
| Assign sample to node $i$ |

**P(assigning a global data point to a local site $i$) = $w_i$**

# Degree Partitioning

Each data point from global dataset is assigned to a local site with probability proportional to the nodes' degree.

**Degree Partitioning**

For each node $i$

Compute the normalized degree $d_i$

Randomly Sample $d_i |Global\ Set|$ points

Assign sample to node $i$

$$d_i = \frac{Degree\ (Node\ i)}{\sum_{j \in Nodes} Degree(Node\ j)}$$

**P(assigning a global data point to a local site $i$) = $d_i$**

# **Results**

# Inferences

• For all the datasets k-means cost ratio was plotted against communication cost.

• k-means cost ratio is defined as the ratio obtained by running Lloyd's algorithm on the coreset and the global data respectively to get two solutions, and computing the ratio between the costs of the two solutions over the global data.

• For all the datasets the k-means cost ratio reduced when communication cost increases. This implies that when larger coresets are generated the solution of the proposed algorithm converges to the k-means solution.

• Same results were obtained on Random and Preferential graph with three types of partitioning namely, uniform partitioning, weighted partitioning and degree partitioning.

• Results were also consistent on different dataset.

# Datasets Used

Following 4 datasets were used

1)  **Corel Image Features:** This dataset contains image features extracted from a Corel There are total 68040 observations.

2) **Spambase Data Set:** There are total 4601 observations with 58 dimensions. There were 1813 marked as spam and 2788 marked as non-spam.

3) **Letter Recognition Dataset:** The dataset was made to identify each of a large number of black-and-white rectangular pixel displays as one of the 26 capital letters in the English alphabet. There are total of 20000 observations each observation has 16 attributes.

4)  **SFpark dataset:** The dataset contains data from three hundred parking lots in San Francisco city fitted with Smart Meters. In the processed data the rows represents Parking lots and columns represents time. There were total 300 rows (one for each parking lot) and 2880 columns (24 hours x 15 time slots x 30 days).

# LETTER DATASET



Random Graph
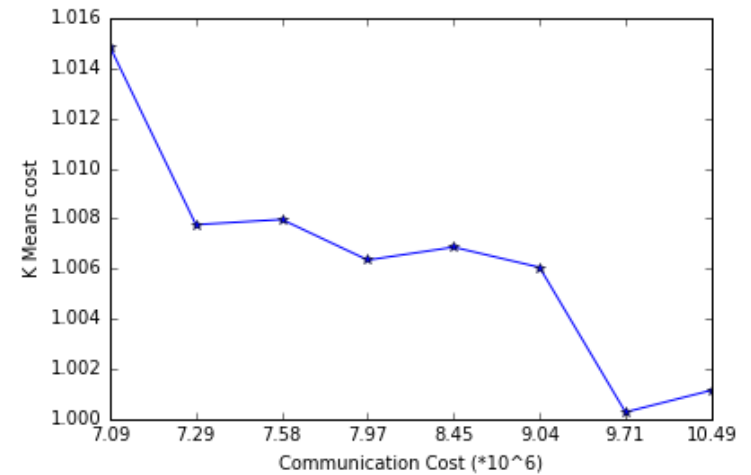Degree Partitioning



Random Graph
Weighted Partitioning



Random Graph
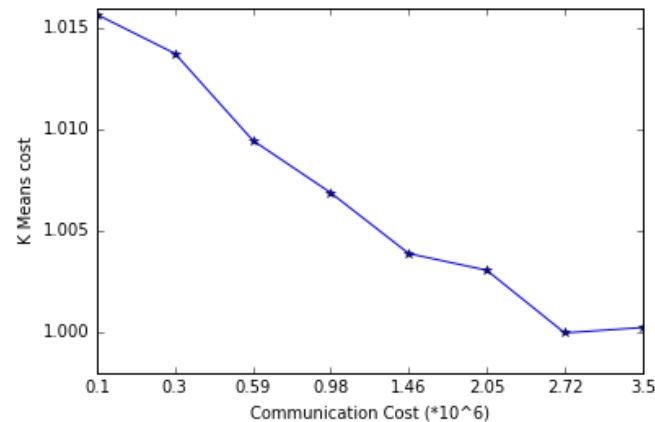Uniform Partitioning

# LETTER DATASET



Preferential Graph
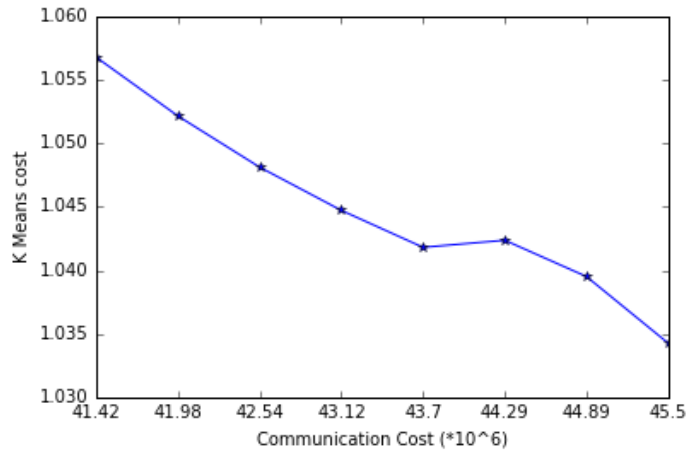Weighted Partitioning

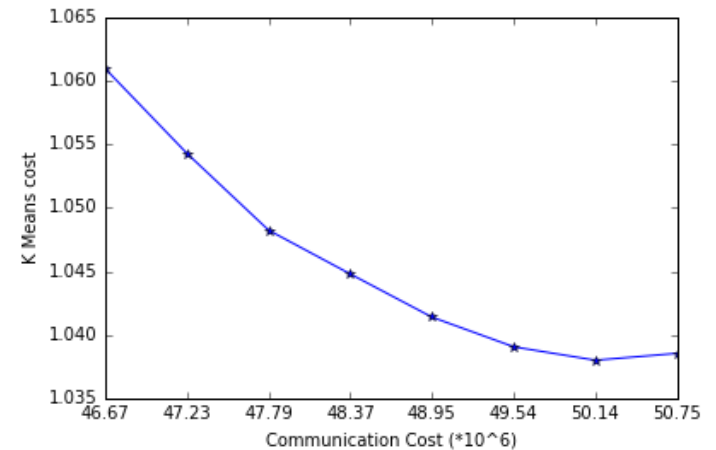Preferential Graph
Uniform Partitioning
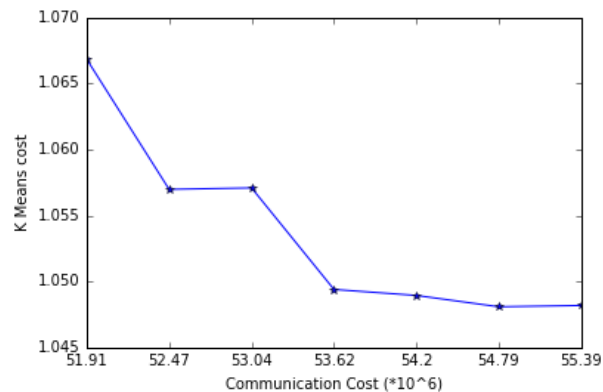
Preferential Graph
Degree Partitioning

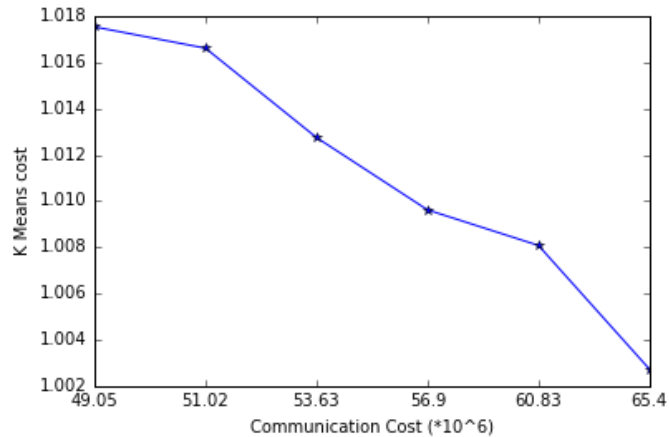# SPAM DATASET



Preferential Graph
Degree Partitioning
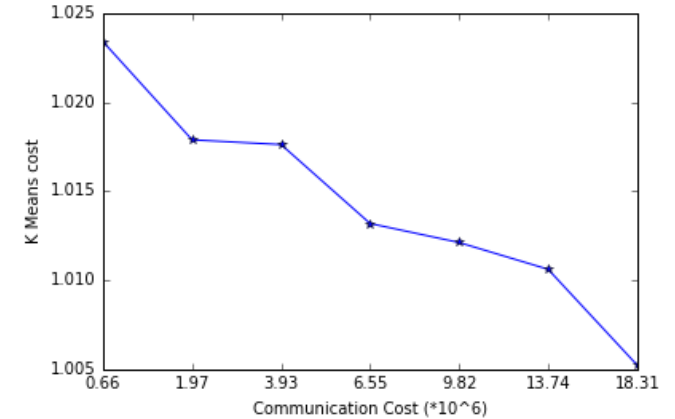
Preferential Graph
Weighted Partitioning

Preferential Graph
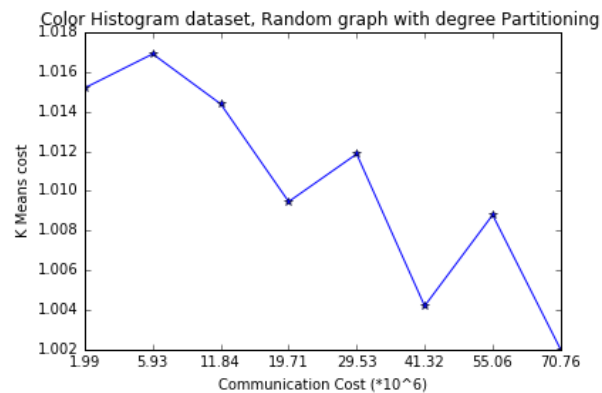Uniform Partitioning

# COLOR HISTOGRAM DATASET
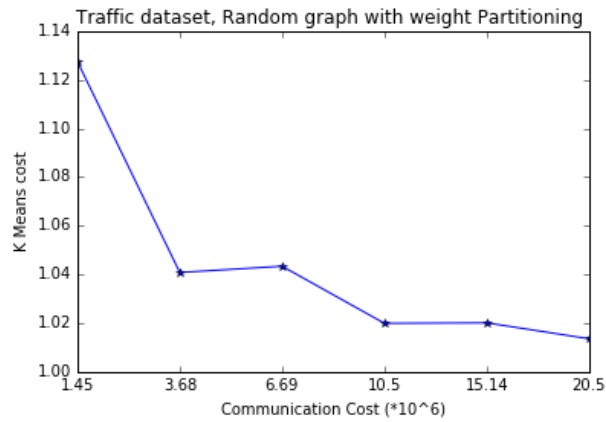


Preferential Graph
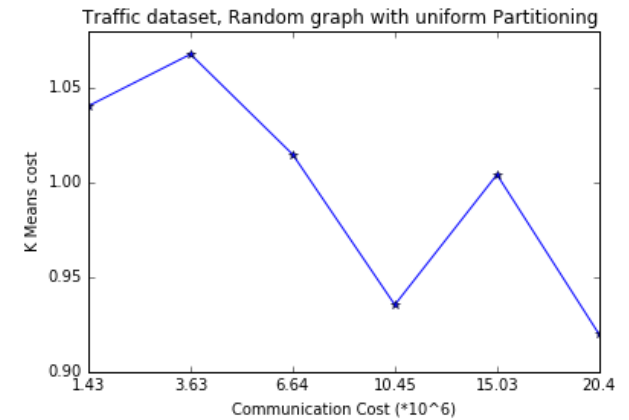Degree Partitioning



Preferential Graph
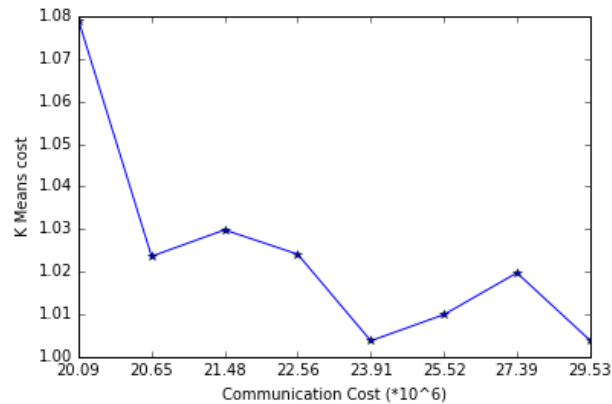Weighted Partitioning



Random Graph
Degree Partitioning

# TRAFFIC DATASET



Random Graph
Weighted Partitioning

Random Graph
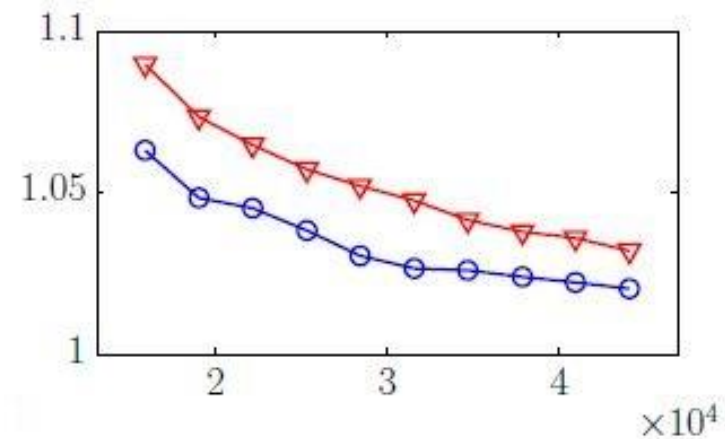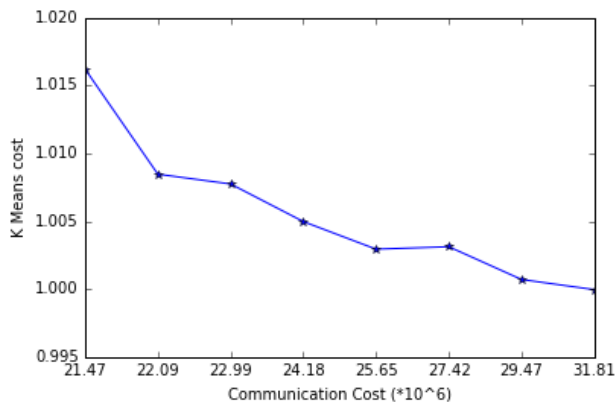Uniform Partitioning

Preferential Graph
Uniform Partitioning

Industrial and Enterprise Systems Engineering

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Summary

- Datasets are getting larger in size, so they are being stored on distributed systems. Increasing the efficiency of the algorithm which work on such distributed systems has gained importance.
- Based on the algorithm proposed in the paper a python code was written to perform distributed k-means clustering.
- 4 different datasets from two different sources was used to evaluate the performance of the algorithm.
- The k-means cost ratio reduces with increase in communication cost.
- We were able to replicate results given in the paper.
- Applications and connections to course

Comparison between two graphs
(Letter Random Weighted)

# References:

[1] Q. Zhang, J. Liu, and W. Wang. Approximate clustering on distributed data streams. In Proceedings of the IEEE International Conference on Data Engineering, 2008.

[2] D. Feldman, A. Sugaya, and D. Rus. An effective coreset compression algorithm for large scale sensor networks. In Proceedings of the International Conference on Information Processing in Sensor Networks, 2012.

[3] K. Bache and M. Lichman. UCI machine learning repository, 2013.

[4] S.R.Boselin Prabhu, S.Sophia, S.Maheswaran, M.Navaneethakrishnan. Real-World Applications of Distributed Clustering Mechanism in Dense Wireless Sensor Networks, 2013.

[5] C. Olston, J. Jiang, and J. Widom. Adaptive filters for continuous queries over distributed data streams. In Proceedings of the ACM SIGMOD International Conference on Management of Data, 2003.

# THE END