



Coreset and sampling approaches for the analysis of very large data sets

Christian Sohler



Big Data - Examples



Image: Wikipedia

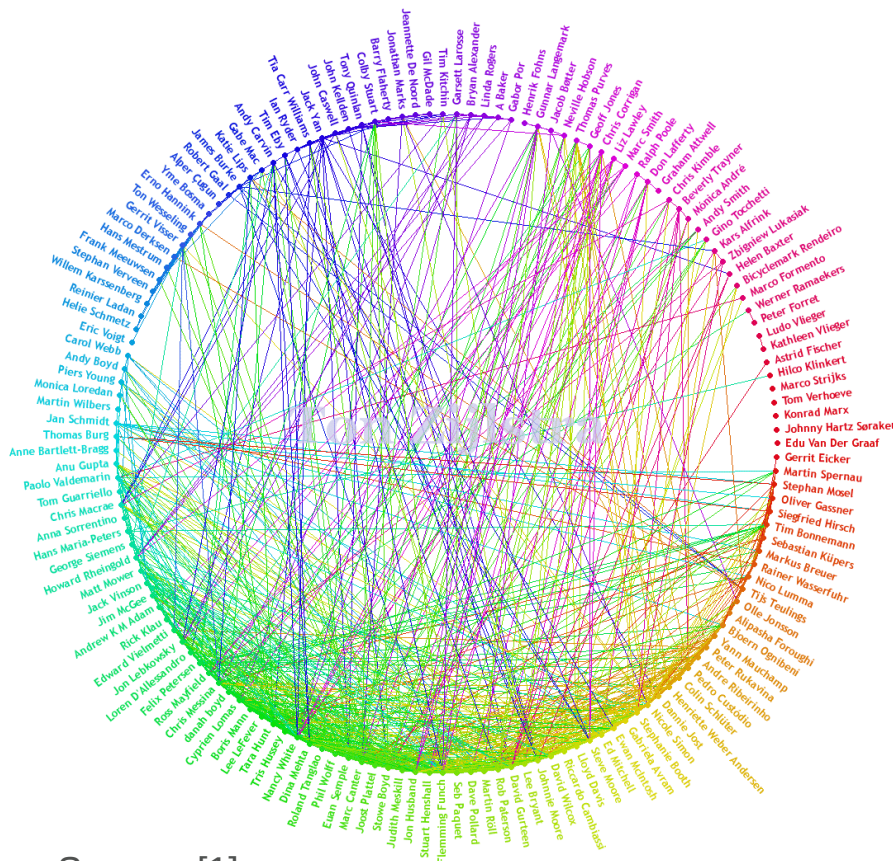
IceCube

- World's biggest neutrino observatory
- Understanding of basic concepts in astrophysics
- Supernovas
- Dark matter
- Black holes

Data size

- approx. 10 TeraByte per data set

Very Large Networks



Social Networks

- Reflect social structures in detail
- Example question: Can we distinguish democratic countries from totalitarian ones by looking at their Facebook structure?

Data size

- GigaByte upto TeraByte (only the graph)
- Data exchange (movies, pictures, etc.) in the Peta-Byte range

Source: [1]

Big Data - Examples



Image: Danard Vincente

The World Wide Web

- Giant source of information
- Very fast access through search engines
- Automatic scoring of the quality of web pages

Data size

- 500 ExaByte in 2009
(Richard Wray, The Guardian)

Big Data Analytics



Image: Roger Smith

Primary Requirements

- Sequential processing
- Sublinear memory requirement
- Sublinear, linear or near linear running time

Secondary Requirements

- Energy efficiency
- Easy to distribute/parallelize
- Limited communication

Outline

Part I: Unstructured Data

- Clustering
- Definition of coresets
- Streaming & distributed algorithms via coresets
- Coreset constructions

Part II: Semi-Structured Data

- Property Testing
- An introductory example
- Local vs. global structure of planar graphs

Cluster analysis

Clustering

- Partition a set of objects into groups such that (ideally)
 - (a) Objects in the same group are similar to each other
 - (b) Objects in different groups are dissimilar to each other

Representation of objects

- Objects are represented by real valued vectors / sets of points

k-Means Clustering

Definition of k-Means Clustering

- Input: Set P of n d-dimensional points
- Output: Set C of k points (centers) that minimize

$$\text{cost}(P,C) = \sum_{p \in P} \min_{c \in C} \|p-c\|^2$$

- The corresponding clustering can be obtained by assigning each point to its nearest center

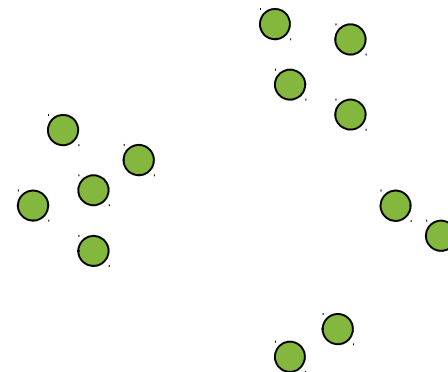
k-Means Clustering

Definition of k-Means Clustering

- Input: Set P of n d-dimensional points
- Output: Set C of k points (centers) that minimize

$$\text{cost}(P,C) = \sum_{p \in P} \min_{c \in C} \|p-c\|^2$$

- The corresponding clustering can be obtained by assigning each point to its nearest center



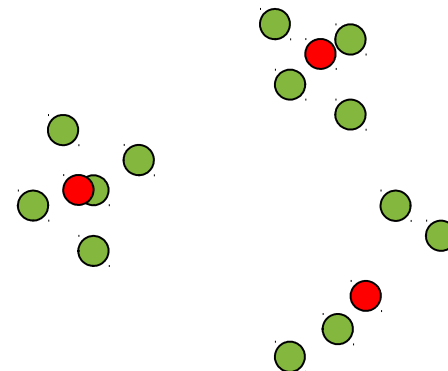
k-Means Clustering

Definition of k-Means Clustering

- Input: Set P of n d-dimensional points
- Output: Set C of k points (centers) that minimize

$$\text{cost}(P,C) = \sum_{p \in P} \min_{c \in C} \|p-c\|^2$$

- The corresponding clustering can be obtained by assigning each point to its nearest center



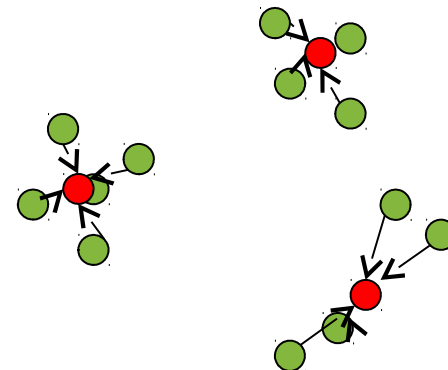
k-Means Clustering

Definition of k-Means Clustering

- Input: Set P of n d-dimensional points
- Output: Set C of k points (centers) that minimize

$$\text{cost}(P,C) = \sum_{p \in P} \min_{c \in C} \|p-c\|^2$$

- The corresponding clustering can be obtained by assigning each point to its nearest center



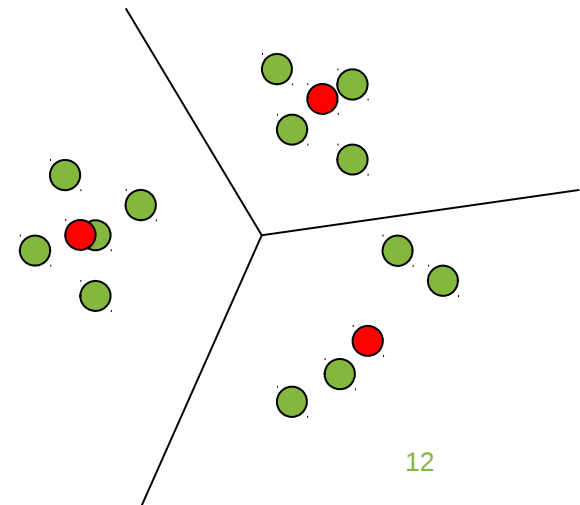
k-Means Clustering

Definition of k-Means Clustering

- Input: Set P of n d-dimensional points
- Output: Set C of k points (centers) that minimize

$$\text{cost}(P,C) = \sum_{p \in P} \min_{c \in C} \|p-c\|^2$$

- The corresponding clustering can be obtained by assigning each point to its nearest center



Coresets – the general idea

Coreset

- A coreset is an approximation of the data with respect to an optimization problem such that a „good“ solution on the coreset is almost as good as a „good“ solution on the original point set
- There is no common definition
- Many approaches can be viewed as coreset approaches (sampling, etc.)

Coresets – a definition

Coresets for k-Means Clustering [Har-Peled, Mazumdar, 2004]

A small weighted set S is called (k, ε) -coreset for P , if for every set C of k centers we have

$$(1-\varepsilon) \text{cost}(P,C) \leq \text{cost}(S,C) \leq (1+\varepsilon) \text{cost}(P,C)$$

In the above definition, the weight of a point is interpreted as its multiplicity.

Coresets – a definition

Coresets for k-Means Clustering [Har-Peled, Mazumdar, 2004]

A small weighted set S is called (k, ε) -coreset for P , if for every set C of k centers we have

$$(1-\varepsilon) \text{cost}(P,C) \leq \text{cost}(S,C) \leq (1+\varepsilon) \text{cost}(P,C)$$

In the above definition, the weight of a point is interpreted as its multiplicity.

Coresets – a definition

Coresets for k-Means Clustering [Har-Peled, Mazumdar, 2004]

A small weighted set S is called (k, ε) -coreset for P , if for every set C of k centers we have

$$(1-\varepsilon) \text{cost}(P,C) \leq \text{cost}(S,C) \leq (1+\varepsilon) \text{cost}(P,C)$$

In the above definition, the weight of a point is interpreted as its multiplicity.

Properties of coresets

- Fast answering of cost queries
- Optimal solution on the coreset yields a $(1+O(\varepsilon))$ -approximation
- α -approximation on the coreset yields an $\alpha \cdot (1+O(\varepsilon))$ -approximation

First Technique: Bounded Movement of Points

Lemma

- Let Opt be the cost of an optimal k -means solution
- If we move a set of points P by an overall l_2^2 -distance of $O(\varepsilon^2 \text{Opt})$ to obtain Q then we have for every set of k centers \hat{C} :

$$(1-\varepsilon) \text{cost}(P, \hat{C}) \leq \text{cost}(Q, \hat{C}) \leq (1+\varepsilon) \text{cost}(P, \hat{C})$$

Main idea [Har-Peled, Mazumdar, 2004]

- Move points such that many are in the same position
- Replace multiple points by weighted points

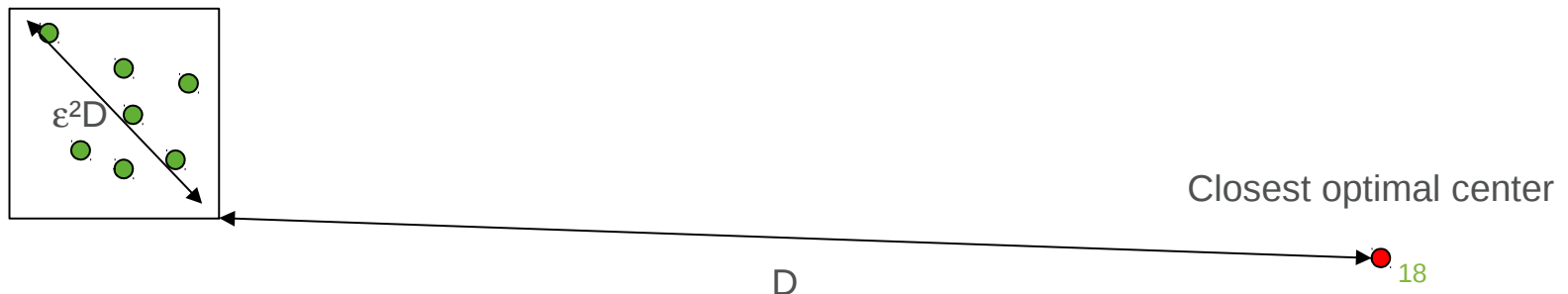
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Mazumdar, 2004]

- There is a coreset of size $O(k \log n/\epsilon d)$ for the k -means problem.

Proof idea:

- Compute a partition of the space into squares s.t. the diagonal of a square is at most ϵ^2 times its l_2^2 -distance to the nearest center of an optimal solution
- Move all points to the center of the square
- Since every point in a cell with l_2^2 -distance D contributes at least D to the optimal solution, we have that the overall movement is $\epsilon^2 \text{ Opt} \Rightarrow \text{Coreset!}$



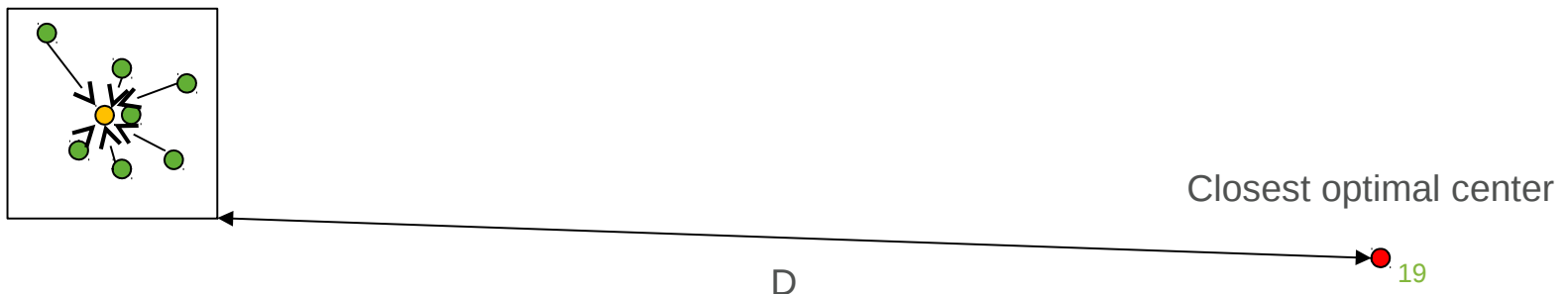
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Mazumdar, 2004]

- There is a coreset of size $O(k \log n / \epsilon d)$ for the k -means problem.

Proof idea:

- Compute a partition of the space into squares s.t. the diagonal of a square is at most ϵ^2 times its l_2^2 -distance to the nearest center of an optimal solution
- Move all points to the center of the square
- Since every point in a cell with l_2^2 -distance D contributes at least D to the optimal solution, we have that the overall movement is $\epsilon^2 \text{ Opt} \Rightarrow \text{Coreset!}$



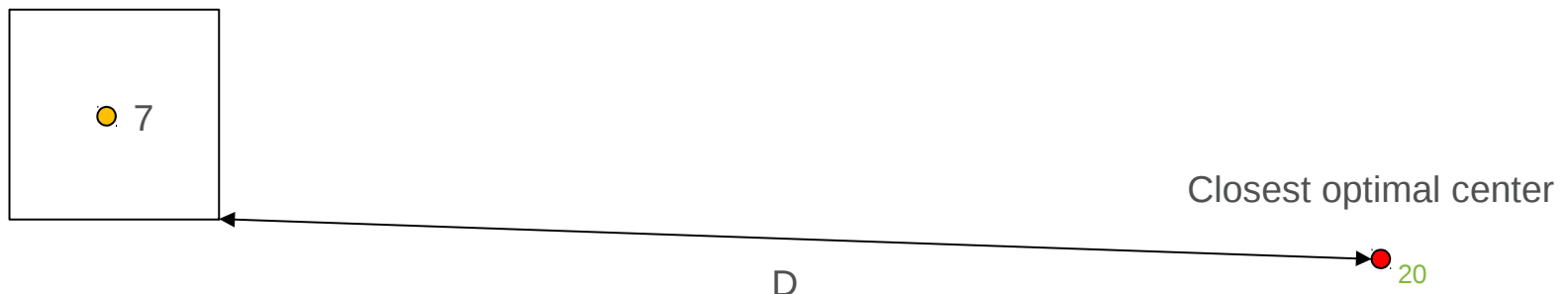
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Mazumdar, 2004]

- There is a coreset of size $O(k \log n / \epsilon d)$ for the k -means problem.

Proof idea:

- Compute a partition of the space into squares s.t. the diagonal of a square is at most ϵ^2 times its l_2^2 -distance to the nearest center of an optimal solution
- Move all points to the center of the square
- Since every point in a cell with l_2^2 -distance D contributes at least D to the optimal solution, we have that the overall movement is $\epsilon^2 \text{ Opt} \Rightarrow \text{Coreset!}$

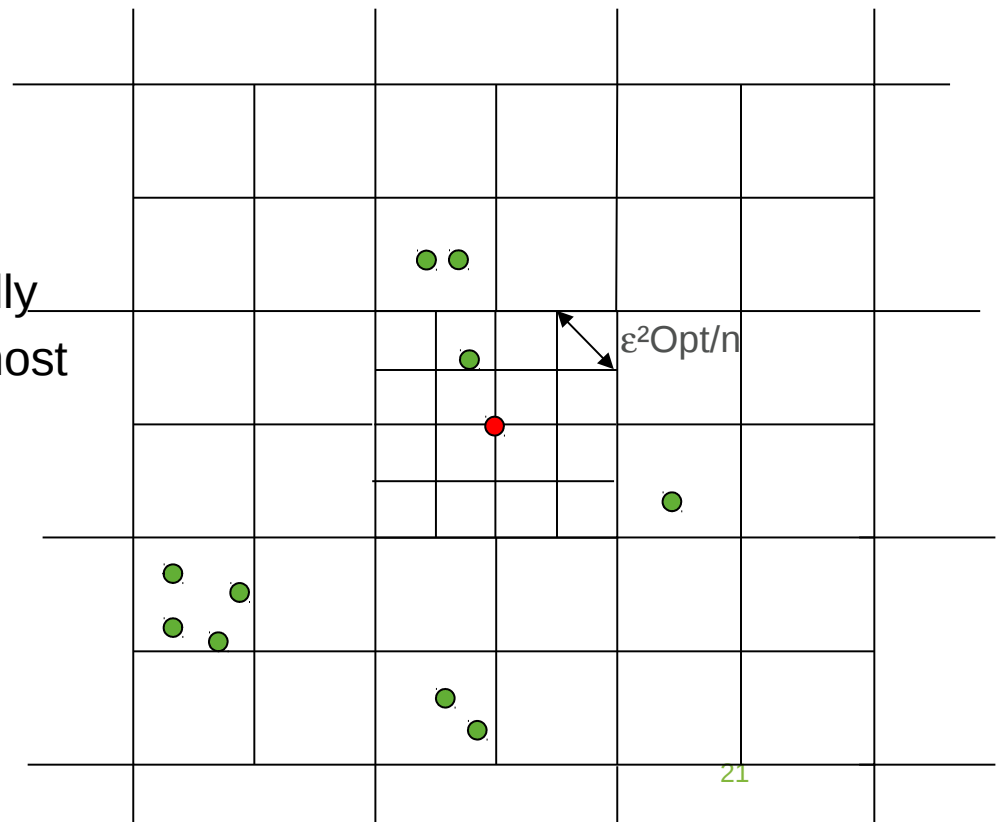


First Technique: Bounded Movement of Points

Theorem [Har-Peled, Mazumdar, 2004]

- There is a coreset of size $O(k \log n / \epsilon d)$ for the k -means problem.

- Smallest cells have diagonal length $\epsilon^2 \text{Opt} / n$
- Cells size increases exponentially
- Furthest point has distance at most Opt
- $\Rightarrow O(k \log n / \epsilon d)$ cells



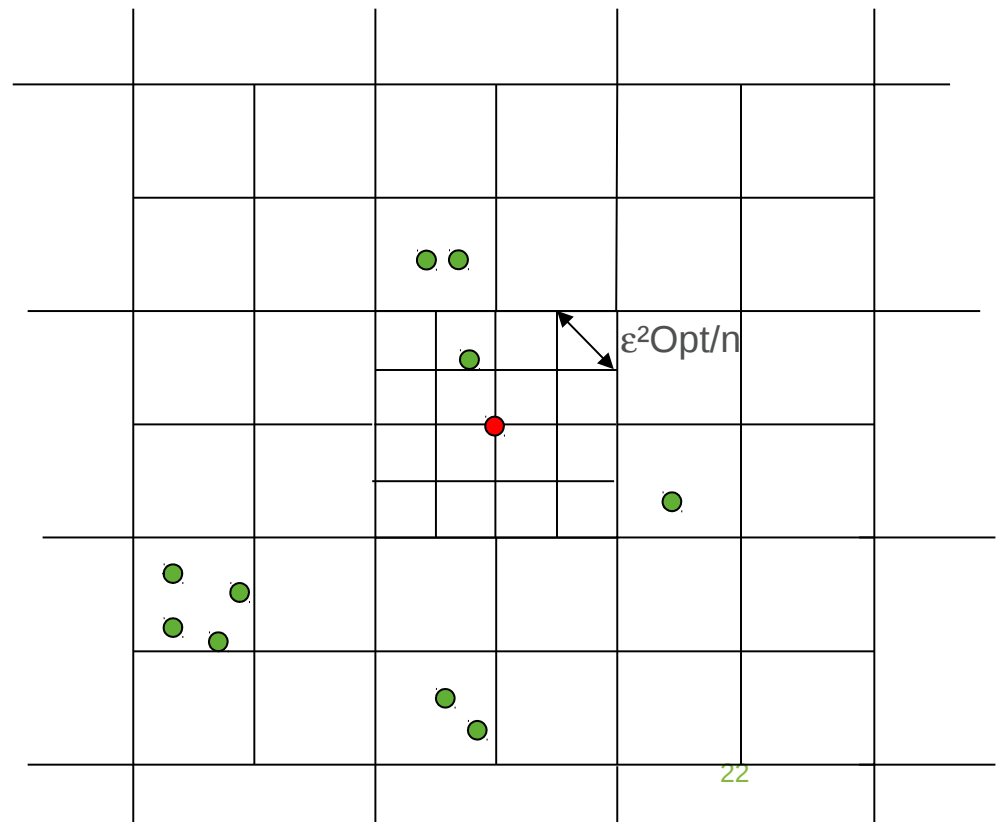
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Mazumdar, 2004]

- There is a coresets of size $O(k \log n/\epsilon d)$ for the k -means problem.

Final observation:

- The construction works, if we replace the optimal solution by an $O(1)$ -approximation



First Technique: Bounded Movement of Points

A different construction [Frahling, S., 2005]

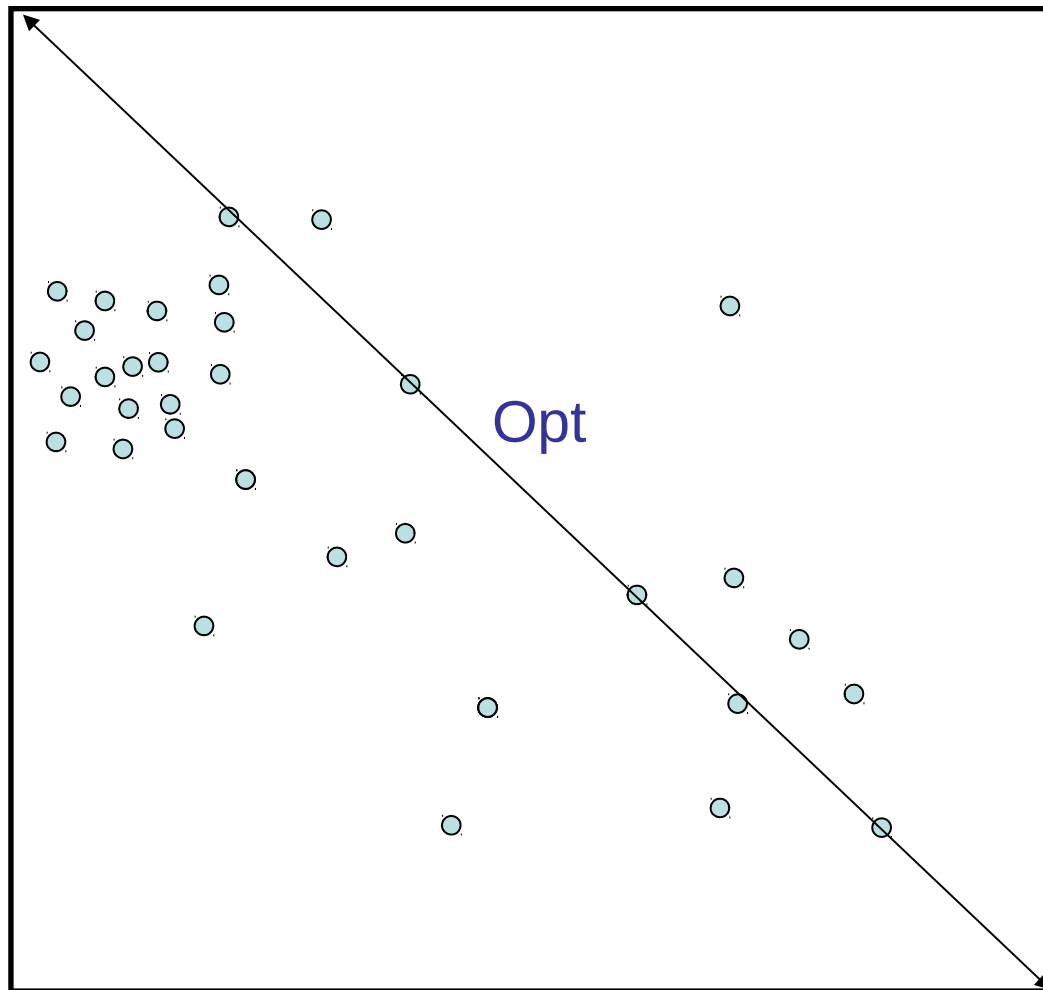
- Idea: Distribute the „movement cost“ equally among all cells
- Movement cost = l_2^2 cell diameter · #points in cell

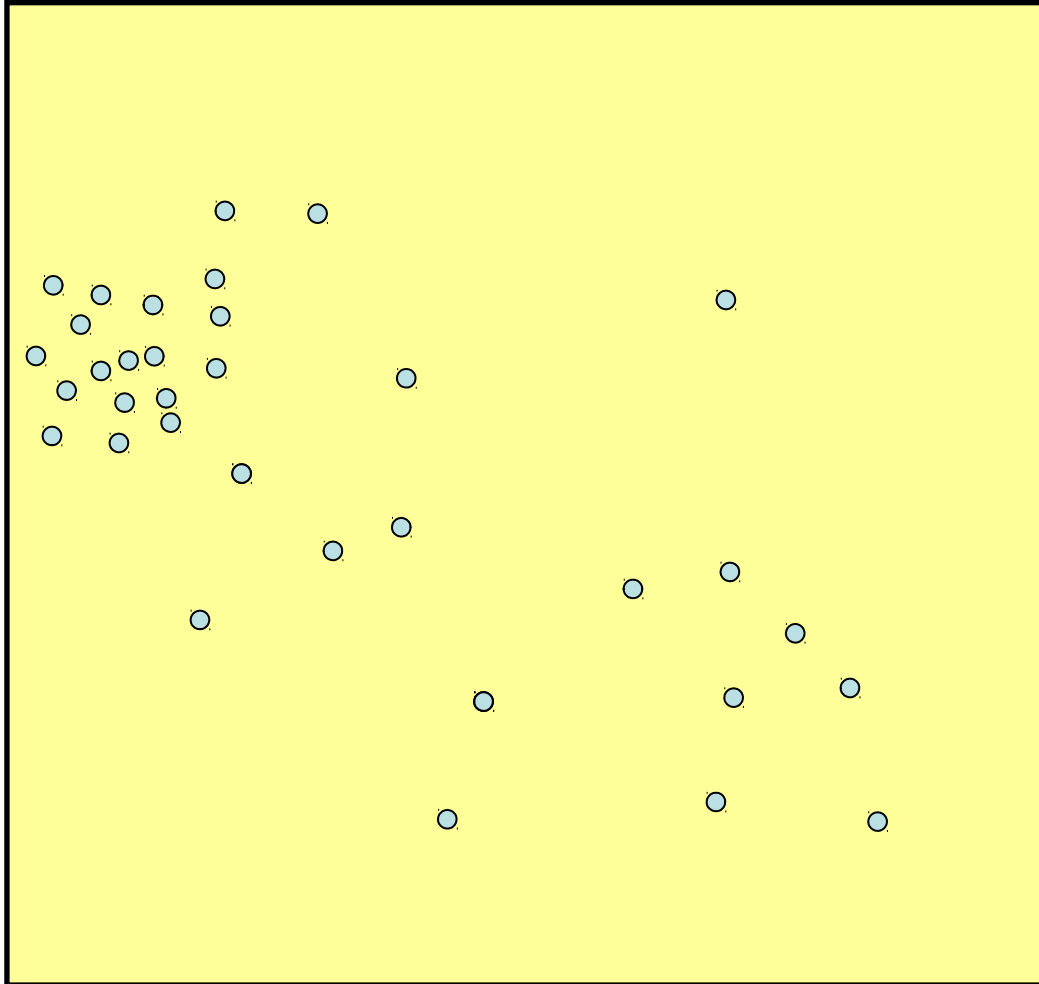
Algorithm

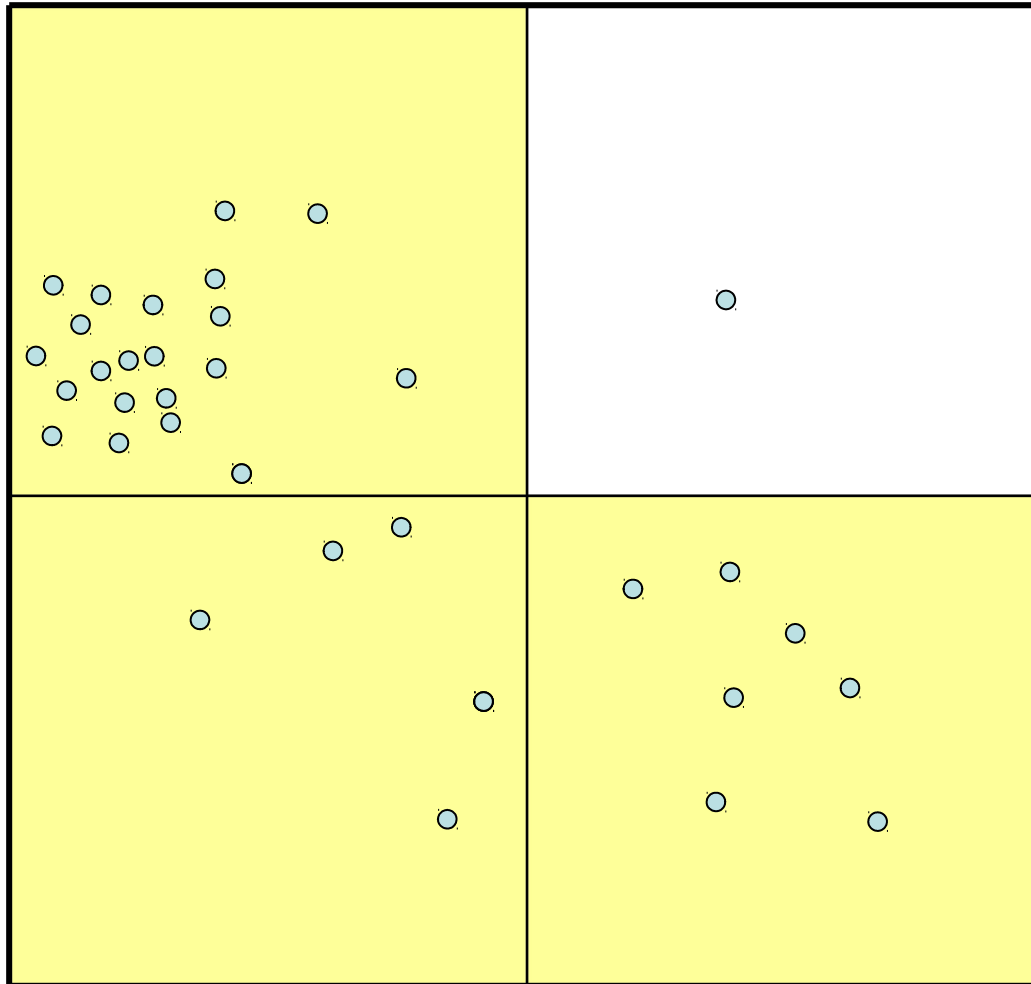
- Guess the cost of an optimal solution, fix appropriate δ
- Call a cell δ -heavy, if l_2^2 cell diameter · #points in cell $> \delta \text{ Opt}$
- Put a coresets point in each minimal δ -heavy cell
- Move each point to a coresets point in the smallest heavy cell it is contained in

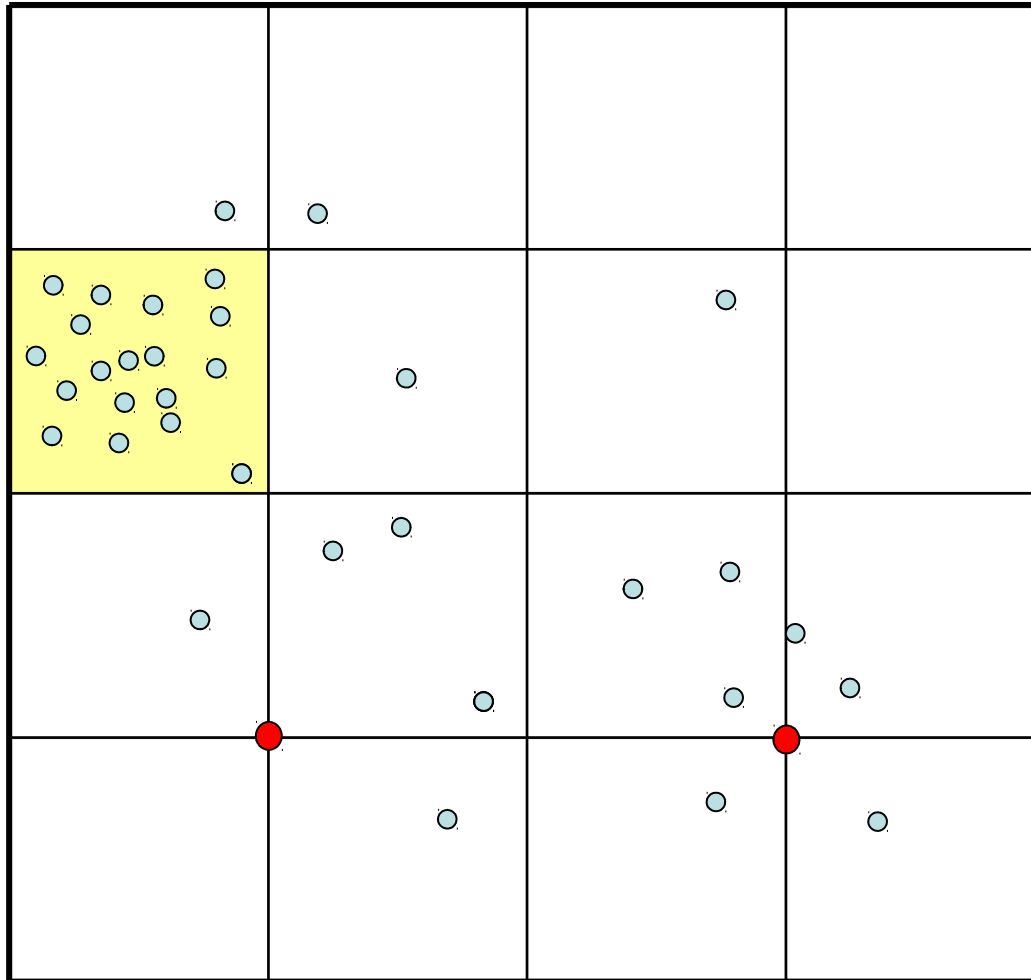
Advantage

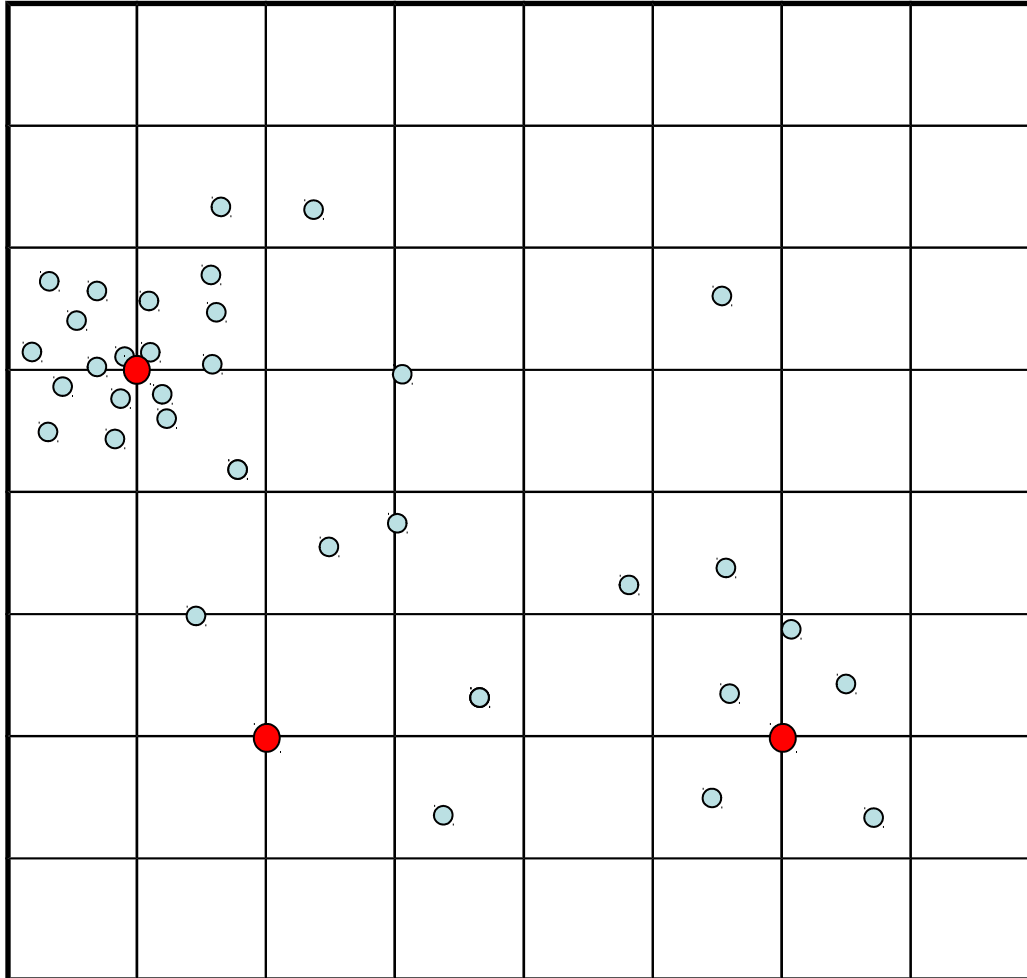
- We do not need to compute an $O(1)$ approximation
- Gives first algorithm for insertion/deletion streams

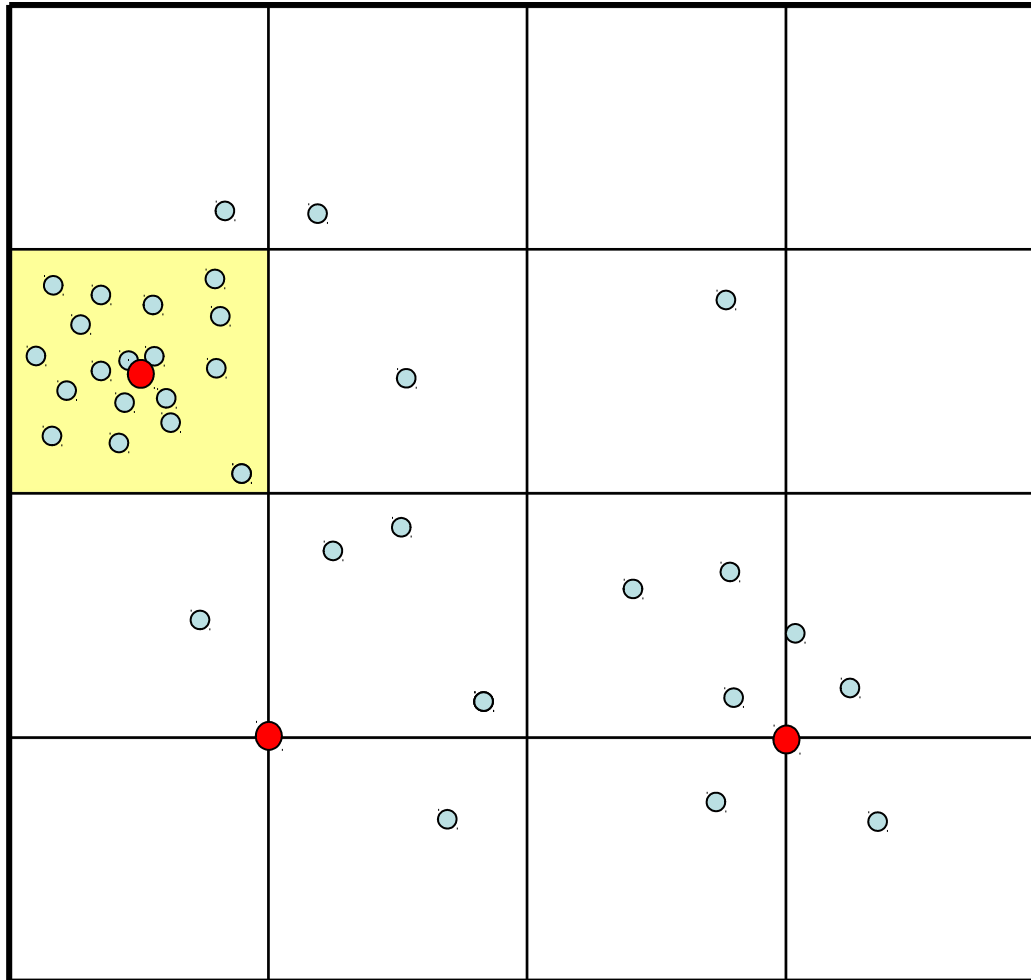


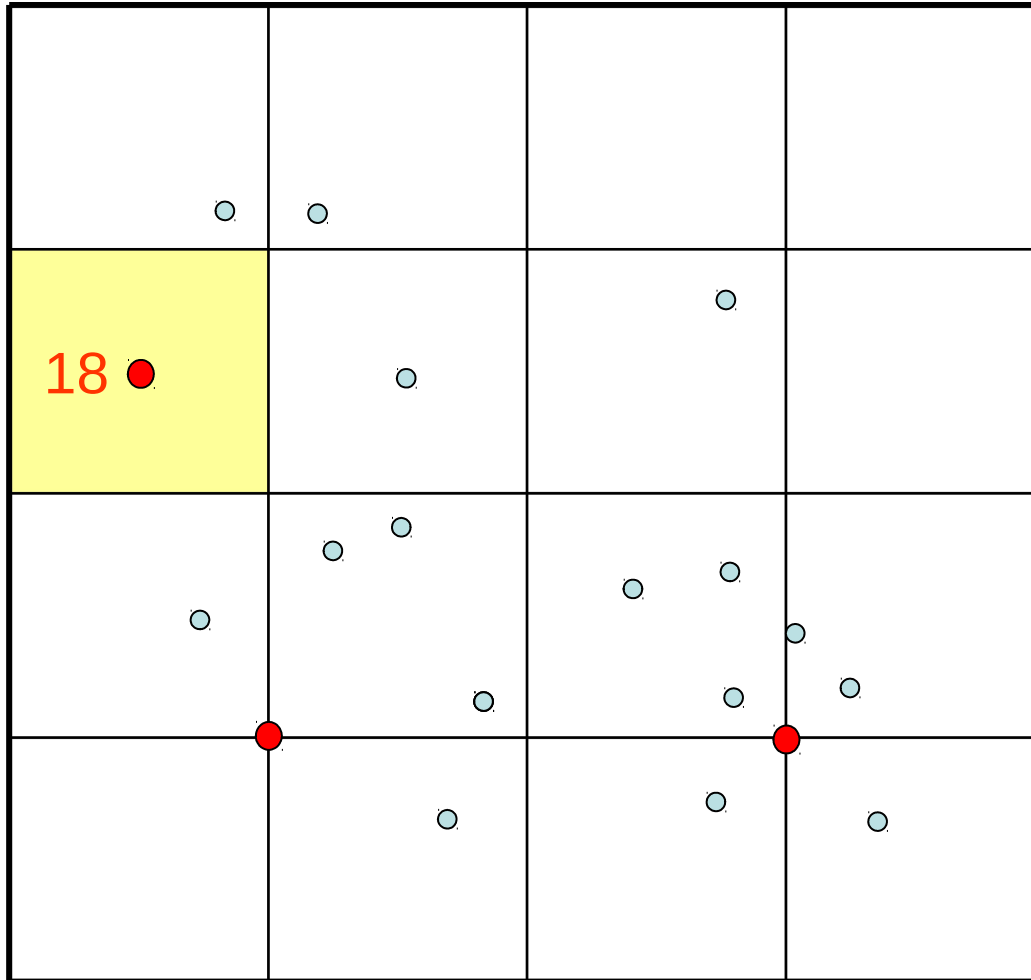


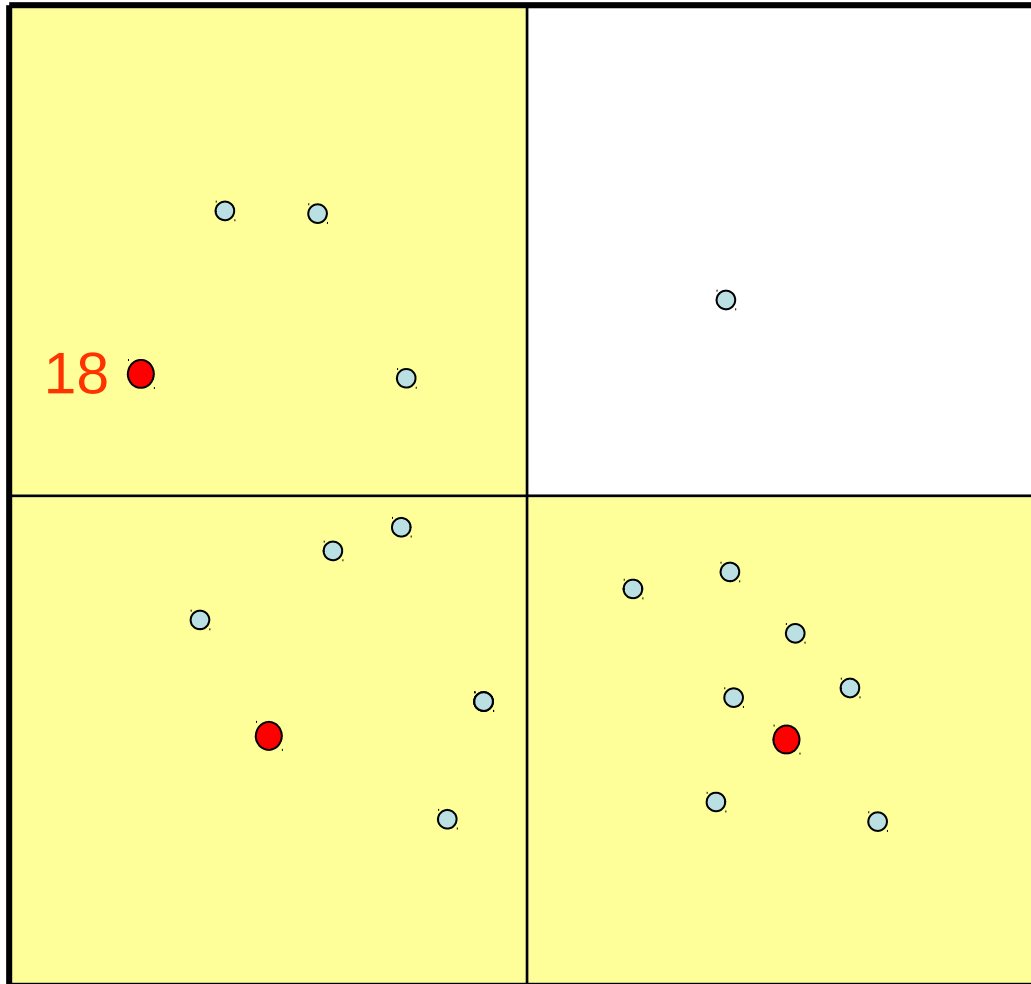


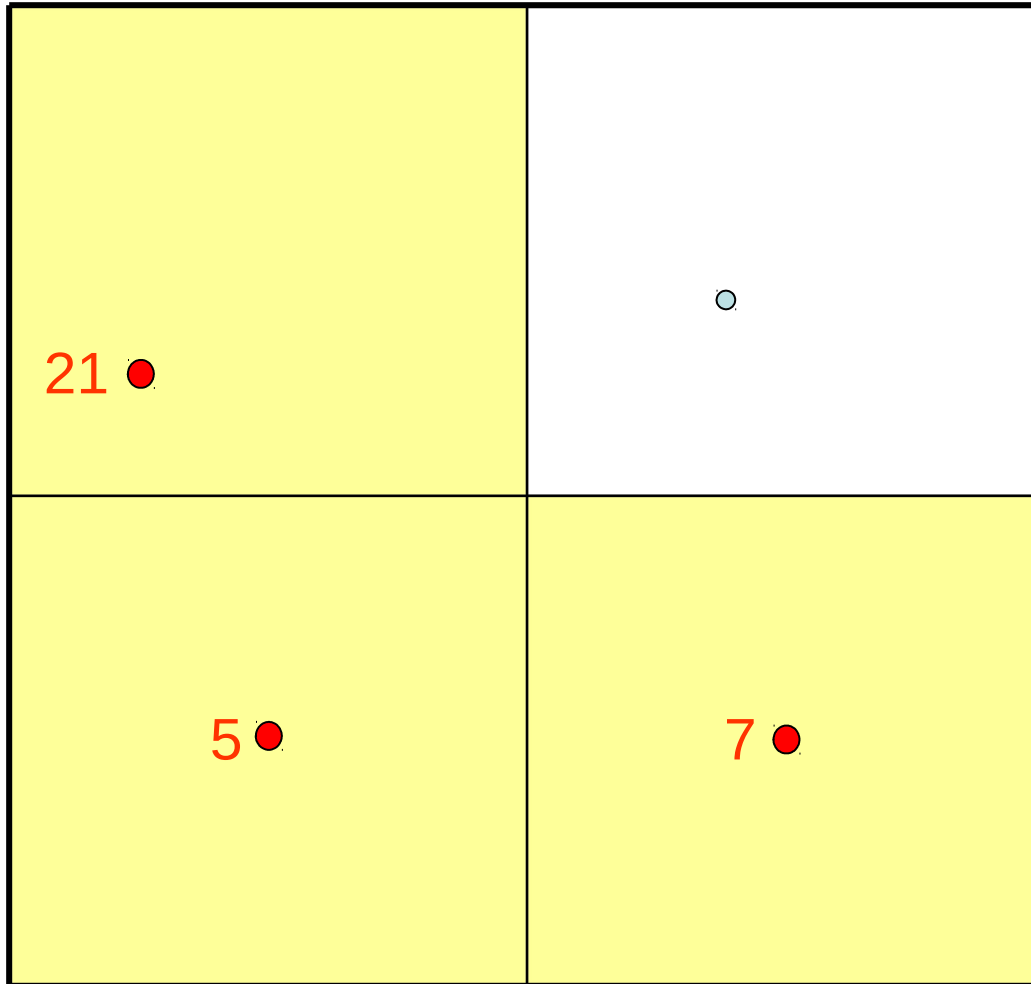


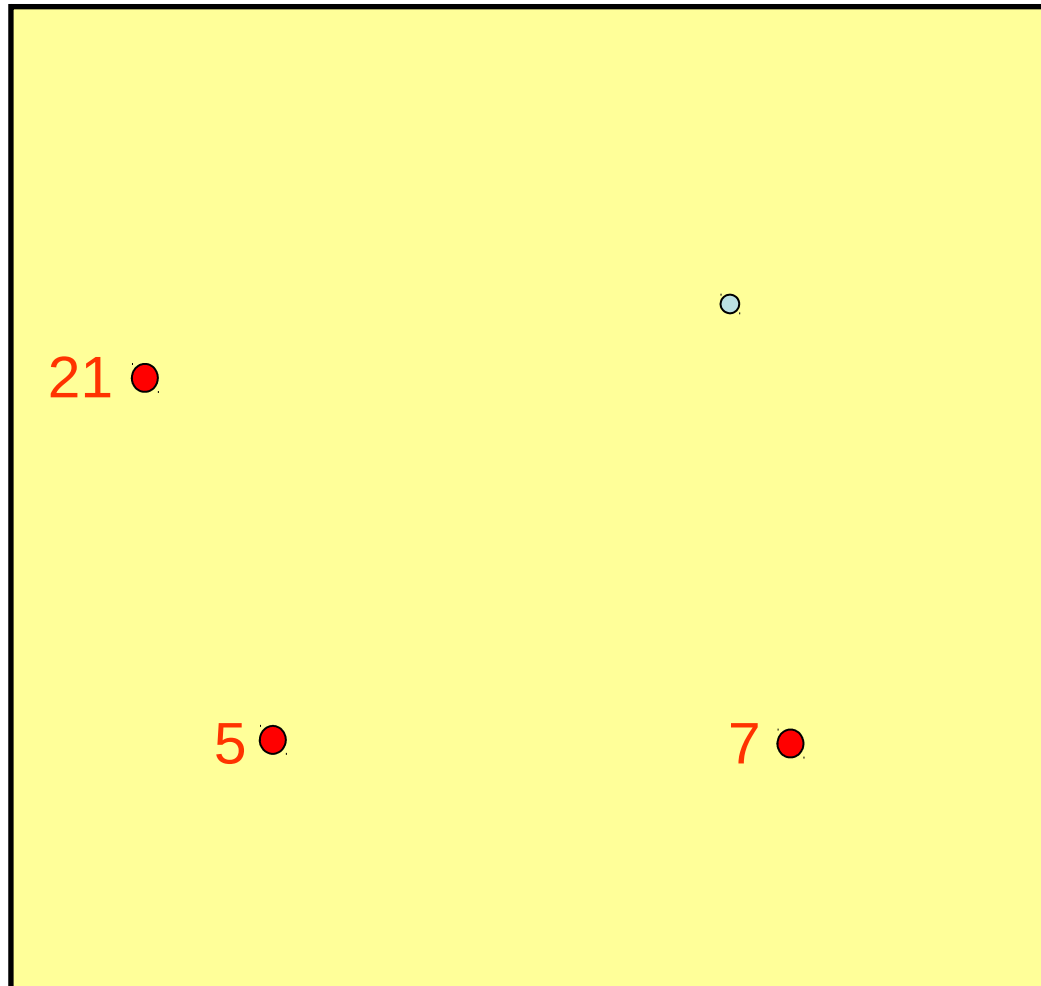


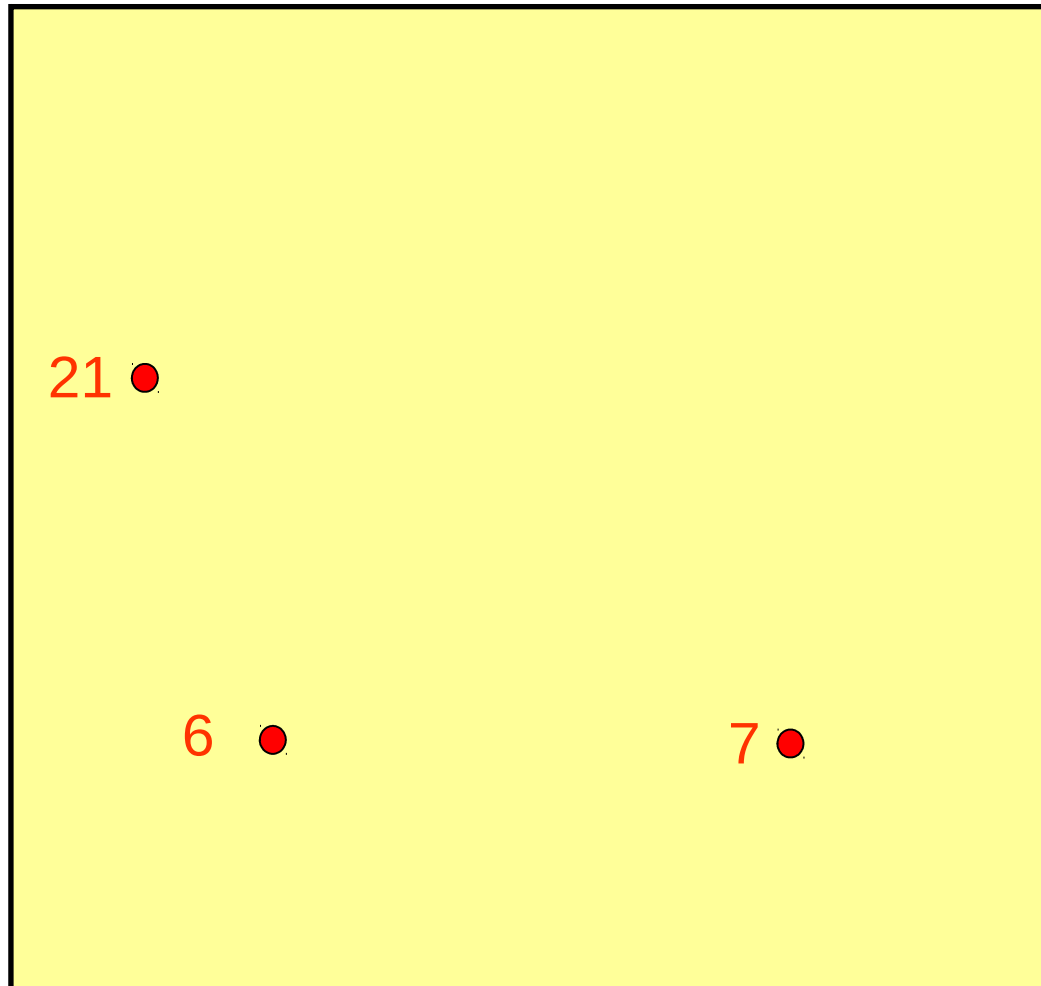












21 ●

6 ●

7 ●

First Technique: Bounded Movement of Points

Theorem [**Frahling, S., 2005**]

- The construction just presented gives a coresets for the k-means problem of size $O(k \log n/\epsilon)$.

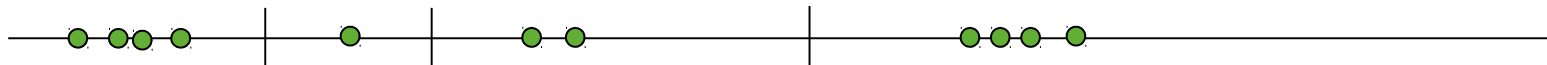
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coresets of size $O(k^3/\epsilon d+1)$ for the k -means problem.

1D Coreset:

- Partition the points into $O(k^2/\epsilon^2)$ intervals such that for each interval the 1-means cost is bounded by $O(\epsilon^2 \text{Opt}/k^2)$



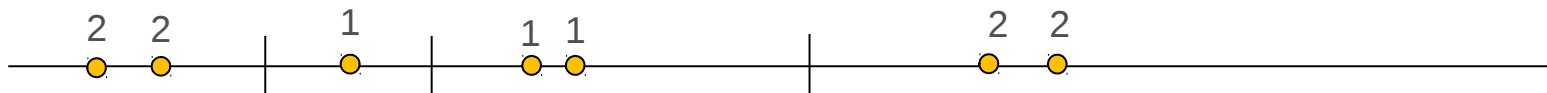
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coresset of size $O(k^3/\epsilon d+1)$ for the k -means problem.

1D Coreset:

- Partition the points into $O(k^2/\epsilon^2)$ intervals such that for each interval the 1-means cost is bounded by $O(\epsilon^2 \text{Opt}/k^2)$
- Replace all points Q in each interval by at most two points S , s.t.



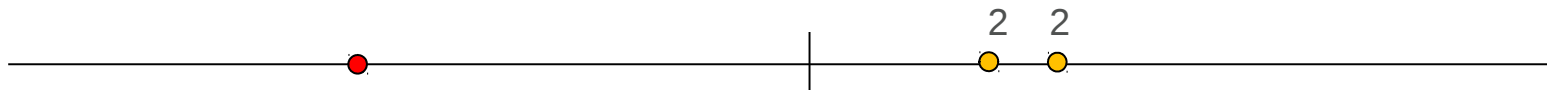
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coreset of size $O(k^3/\epsilon d+1)$ for the k -means problem.

1D Coreset:

- Partition the points into $O(k^2/\epsilon^2)$ intervals such that for each interval the 1-means cost is bounded by $O(\epsilon^2 \text{Opt}/k^2)$
- Replace all points Q in each interval by at most two points S , s.t. if the nearest center is not contained in the interval, then
$$\text{cost}(Q,c) = \text{cost}(S,c)$$



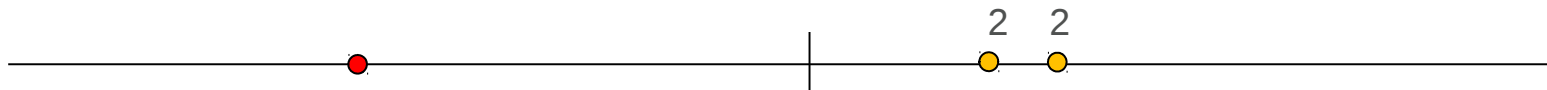
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coreset of size $O(k^3/\epsilon d+1)$ for the k -means problem.

1D Coreset:

- Partition the points into $O(k^2/\epsilon^2)$ intervals such that for each interval the 1-means cost is bounded by $O(\epsilon^2 \text{Opt}/k^2)$
- Replace all points Q in each interval by at most two points S , s.t. if the nearest center is not contained in the interval, then
$$\text{cost}(Q,c) = \text{cost}(S,c)$$



First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coreset of size $O(k^3/\epsilon d+1)$ for the k -means problem.

1D Coreset:

- Partition the points into $O(k^2/\epsilon^2)$ intervals such that for each interval the 1-means cost is bounded by $O(\epsilon^2 \text{Opt}/k^2)$
- Replace all points Q in each interval by at most two points S , s.t. if the nearest center is not contained in the interval, then
$$\text{cost}(Q,c) = \text{cost}(S,c)$$
- Error only within the intervals that contain a center



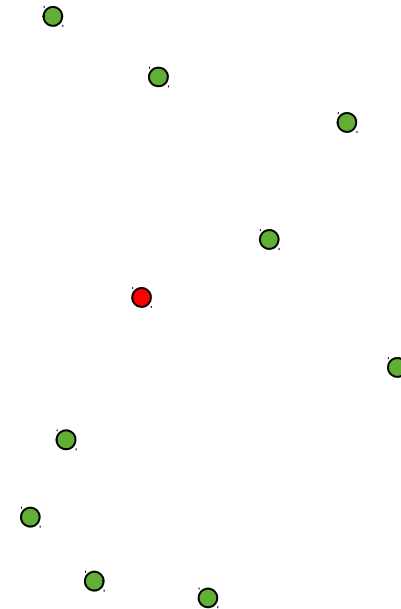
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coreset of size $O(k^3/\epsilon d+1)$ for the k-means problem.

Coresets for higher dimensions:

- Move points to $O(1/\epsilon d-1)$ lines through each center of $O(1)$ -approximation (overall l_2^2 -movement is $\epsilon^2 \text{Opt}$)
- Compute a coreset for each line



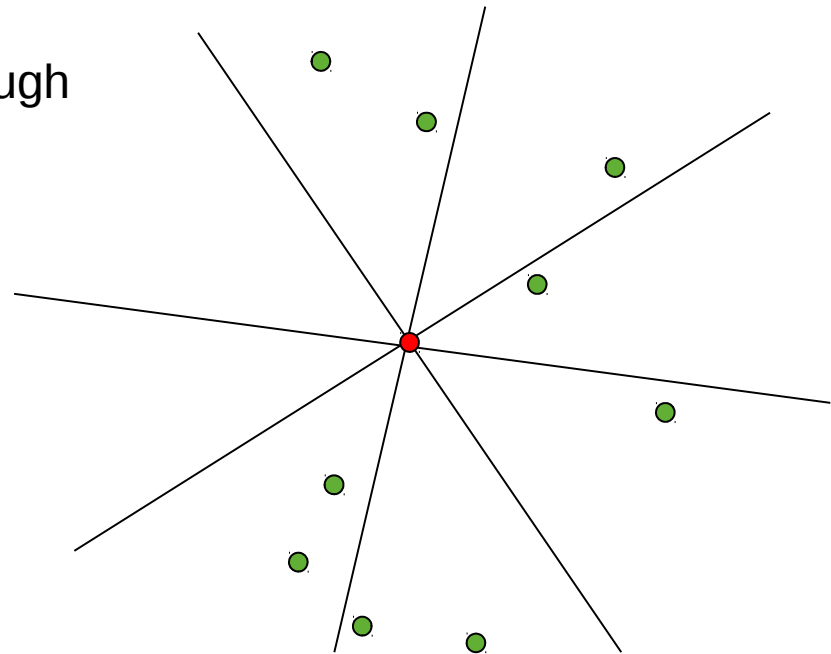
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coreset of size $O(k^3/\epsilon d+1)$ for the k-means problem.

Coresets for higher dimensions:

- Move points to $O(1/\epsilon d-1)$ lines through each center of $O(1)$ -approximation (overall l_2^2 -movement is $\epsilon^2 \text{Opt}$)
- Compute a coreset for each line



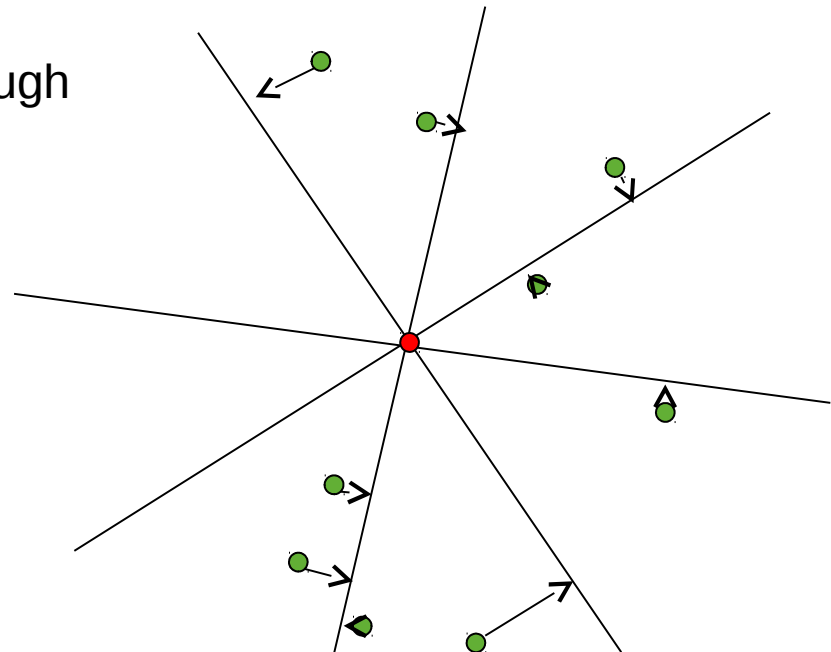
First Technique: Bounded Movement of Points

Theorem [Har-Peled, Kushal, 2005]

- There is a coreset of size $O(k^3/\epsilon d+1)$ for the k-means problem.

Coresets for higher dimensions:

- Move points to $O(1/\epsilon d-1)$ lines through each center of $O(1)$ -approximation (overall l_2^2 -movement is $\epsilon^2 \text{Opt}$)
- Compute a coreset for each line using the 1D construction



Second Technique: Sampling

Lemma

Let $0 < \epsilon < 1$. Let P be a point set in the unit ball. Let C be a fixed set of k centers inside the unit ball. Let S be a sample set taken from P uniformly at random. Then

$$\Pr[|\text{cost}(P,C) - |P|/|S| \text{cost}(S,C)| > \epsilon n] < 2 \exp(- \epsilon^2 |S| / 12)$$

Proof:

- Apply Hoeffding's bound

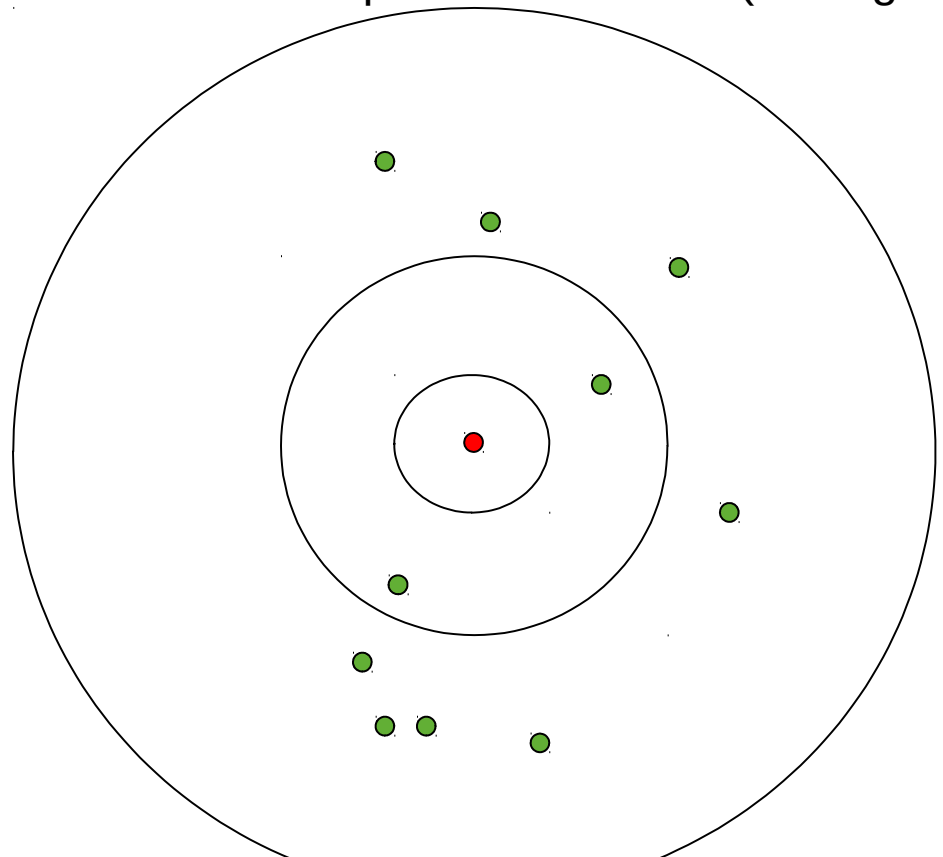
Second Technique: Sampling

Theorem [Chen, 2009]

Let $0 < \epsilon < 1$. There is an (ϵ, k) -coreset for the k -means problem of size $\tilde{O}(dk^2 \log n / \epsilon^2)$.

Algorithm:

- Partition space around $O(1)$ -approximation into exponentially increasing rings
- From each ring take a uniform sample



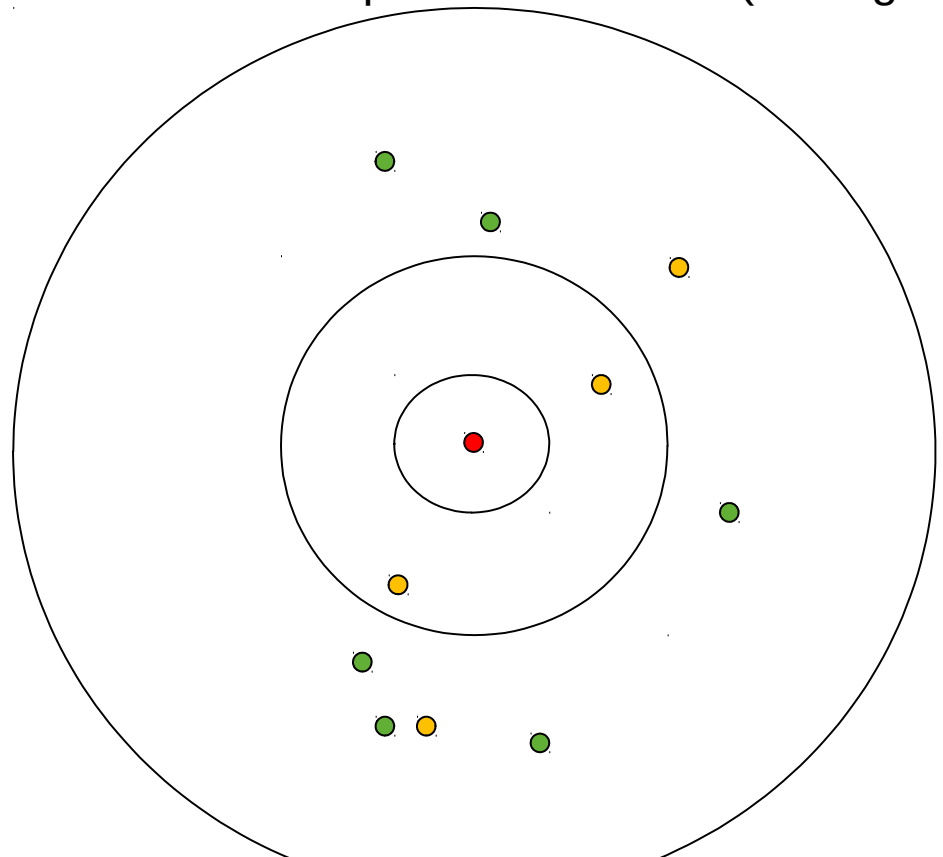
Second Technique: Sampling

Theorem [Chen, 2009]

Let $0 < \epsilon < 1$. There is an (ϵ, k) -coreset for the k -means problem of size $\tilde{O}(dk^2 \log n / \epsilon^2)$.

Algorithm:

- Partition space around $O(1)$ -approximation into exponentially increasing rings
- From each ring take a uniform sample
- Weight the sample appropriately



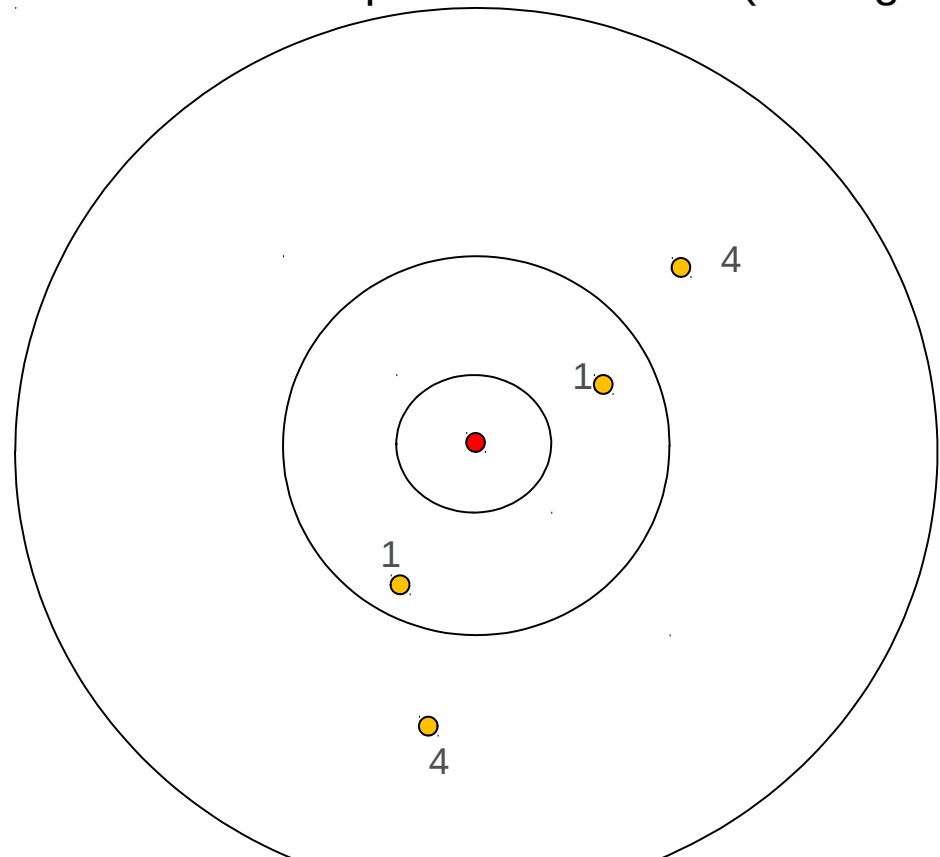
Second Technique: Sampling

Theorem [Chen, 2009]

Let $0 < \epsilon < 1$. There is an (ϵ, k) -coreset for the k -means problem of size $\tilde{O}(dk^2 \log n / \epsilon^2)$.

Algorithm:

- Partition space around $O(1)$ -approximation into exponentially increasing rings
- From each ring take a uniform sample
- Weight the sample appropriately



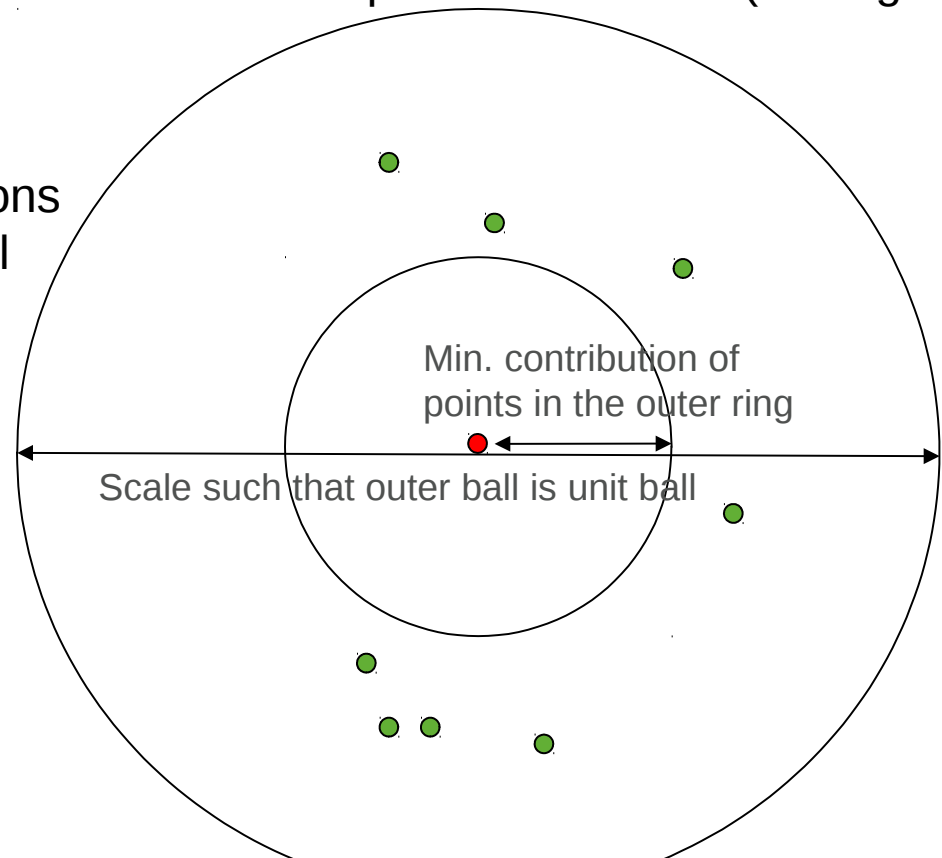
Second Technique: Sampling

Theorem [Chen, 2009]

Let $0 < \epsilon < 1$. There is an (ϵ, k) -coreset for the k -means problem of size $\tilde{O}(dk^2 \log n / \epsilon^2)$.

Proof:

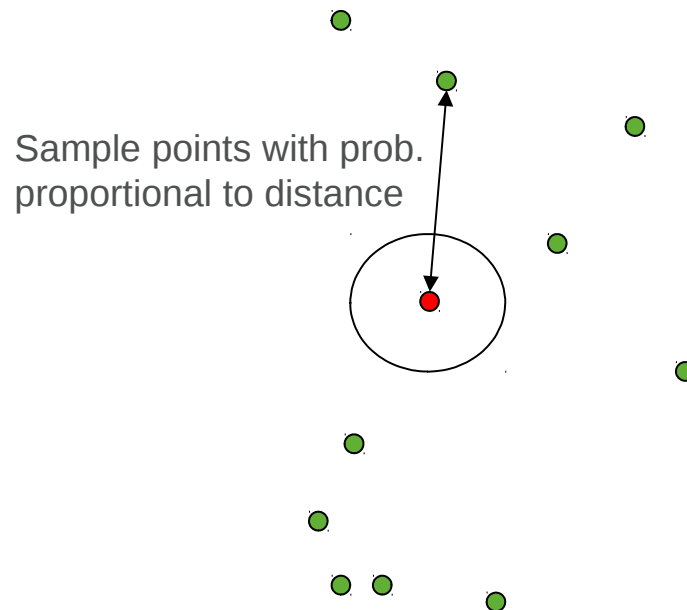
- Define a net of candidate solutions and apply sampling lemma to all solutions from the net
- Charge sampling error to the contribution of the points inside the ring



Technique 2.5: Importance Sampling

Idea [Feldman, Monemizadeh, S., 2005]

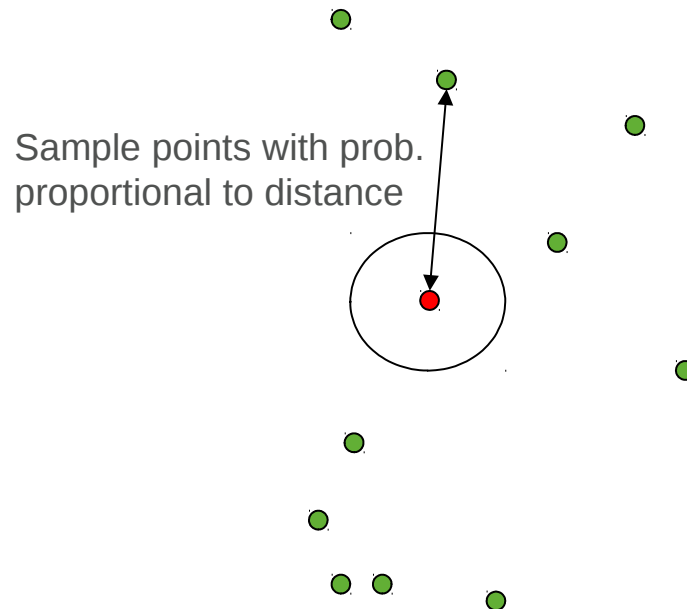
- Sample uniformly the inner ring
- Sample all remaining points according to $\Pr[p] = \text{cost}(\{p\}, C_{\text{Opt}}) / \text{Opt}$ and assign weight $1/\Pr[p]$.



Technique 2.5: Importance Sampling

Idea [Feldman, Monemizadeh, S., 2005]

- Sample uniformly the inner ring
- Sample all remaining points according to $\Pr[p] = \text{cost}(\{p\}, C_{\text{Opt}}) / \text{Opt}$ and assign weight $1/\Pr[p]$.
- Analyze using Hoeffding's bound for all solutions on an appropriate net



Technique 2.5: Importance Sampling

Idea [Feldman, Monemizadeh, S., 2005]

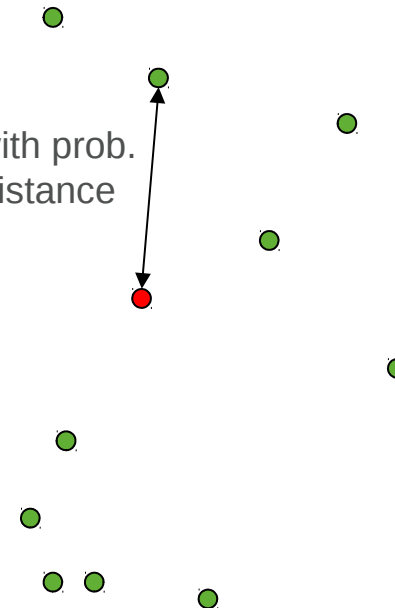
- Sample uniformly the inner ring
- Sample all remaining points according to $\Pr[p] = \text{cost}(\{p\}, C_{\text{Opt}}) / \text{Opt}$ and assign weight $1/\Pr[p]$.
- Analyze using Hoeffding's bound for all solutions on an appropriate net

Sample points with prob.
proportional to distance

Removing the inner ring

[Schulman, Langberg, 2010]

- $\Pr[p] = \text{cost}(\{p\}, C_{\text{Opt}}) / (2 \text{Opt}) + 1/(2n)$



Technique 2.5: Importance Sampling

Definition [Langberg, Schulman, 2010]

The *sensitivity* $\gamma(p)$ of a point $p \in P$ with respect to k -means clustering is defined as

$$\gamma(p) = \sup_{C=\{c_1, \dots, c_k\} \subseteq \mathbb{R}^d} \text{cost}(\{p\}, C) / \text{cost}(P, C).^*$$

The *total sensitivity* $T(P)$ of P is defined as $T(P) = \sum_{p \in P} \gamma(p)$

* assuming that $\text{cost}(\{p\}, C) / \text{cost}(P, C) = 0$, if $\text{cost}(\{p\}, C) = 0$.

Technique 2.5: Importance Sampling

Theorem [**Feldman, Langberg, 2011**]

- If we choose $s \geq d^*/\epsilon^2 \Gamma(P)^2$ points with probability proportional to $\gamma(p)$ then w.p. at least $2/3$ the resulting point set is a (k, ϵ) -coreset.
- Here, d^* is the VC dimension of the range space of the family of k balls.

Third Technique: Reducing the Variance

Coresets with Offset for k-Means [Feldman, Schmidt, S., 2013]

A pair (S, Δ) of a small weighted set S and a value Δ is called *(k, ϵ) -coreset with offset* for P , if for every set C of k centers we have

$$(1-\epsilon) \text{cost}(P, C) \leq \text{cost}(S, C) + \Delta \leq (1+\epsilon) \text{cost}(P, C)$$

In the above definition, the weight of a point is interpreted as its multiplicity.

Third Technique: Reducing the Variance

Coresets with Offset for k-Means [Feldman, Schmidt, S., 2013]

A pair (S, Δ) of a small weighted set S and a value Δ is called *(k, ε) -coreset with offset* for P , if for every set C of k centers we have

$$(1-\varepsilon) \text{cost}(P, C) \leq \text{cost}(S, C) + \Delta \leq (1+\varepsilon) \text{cost}(P, C)$$

In the above definition, the weight of a point is interpreted as its multiplicity.

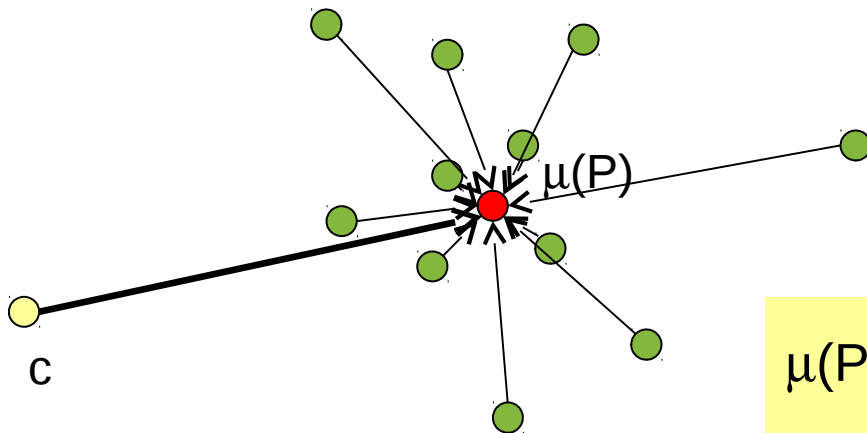
Remarks

- Δ may depend on P , k and ε
- This is **not** an additive approximation
- The variance in S is smaller than in P (by the amount of Δ)

Third Technique: Reducing the Variance

The „Magic Formula“ for k-Means [Zhang, Ramakrishnan, Livny, 96]

$$\sum_{p \in P} \|p - c\|^2 = |P| \cdot \|c - \mu(P)\|^2 + \sum_{p \in P} \|p - \mu(P)\|^2$$

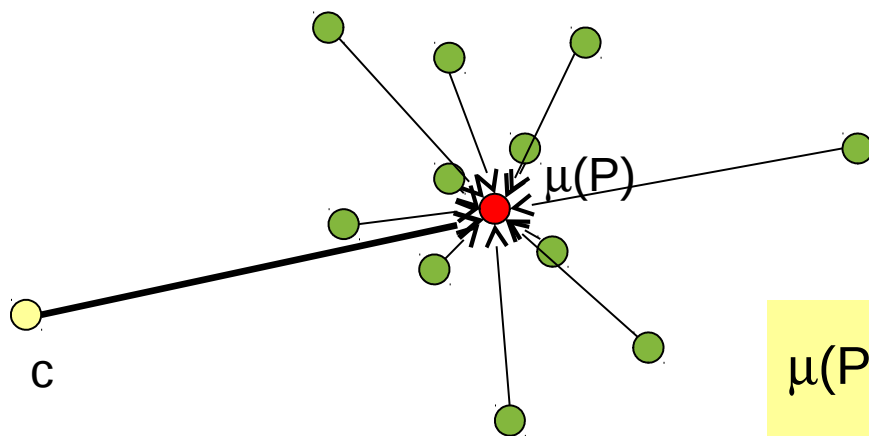


$$\mu(P) = \sum_{p \in P} p / |P|$$

Third Technique: Reducing the Variance

The „Magic Formula“ for k-Means [Zhang, Ramakrishnan, Livny, 96]

$$\sum_{p \in P} \|p - c\|^2 = |P| \cdot \|c - \mu(P)\|^2 + \sum_{p \in P} \|p - \mu(P)\|^2$$



Exact coresets for 1-means

Pair $\mu(P)$ with weight $|P|$ and the value $\Delta = \sum \|p - \mu(P)\|^2$

$$\mu(P) = \sum_{p \in P} p / |P|$$

Third Technique: Reducing the Variance

Definition

- Let $\text{Opt}_k(P)$ denote the optimal clustering cost with k centers. A set P that satisfies $\text{Opt}_1(P) \leq (1+\varepsilon)\text{Opt}_k(P)$ is called *pseudorandom*.

Intuition

- A set is pseudorandom, if there is no good clustering (it behaves like a random point set).
- In particular, if we take an arbitrary partition into k clusters, their weighted means will be close to the mean of P

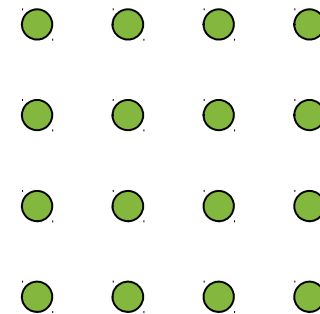
Third Technique: Reducing the Variance

A construction for $k > 1$

- We recursively partition P into subsets $\{P_i\}$ by applying the following rule:
- Let P_i be the set with the highest cost $\text{Opt}_1(P_i)$
- **If** P_i is pseudorandom **then** replace P_i by $\mu(P_i)$
- **else** find an optimal clustering of P_i into k sets $\{Q_j\}$ and replace P_i by the $\{Q_j\}$

$[\sum_j \text{Opt}_1(Q_j)$ decreases by a factor of $1+\epsilon$]

- If $\sum_j \text{Opt}_1(Q_j) \leq \epsilon^2 \text{Opt}_k(P)$ then replace all $\{Q_j\}$ by $\{\mu_j\}$



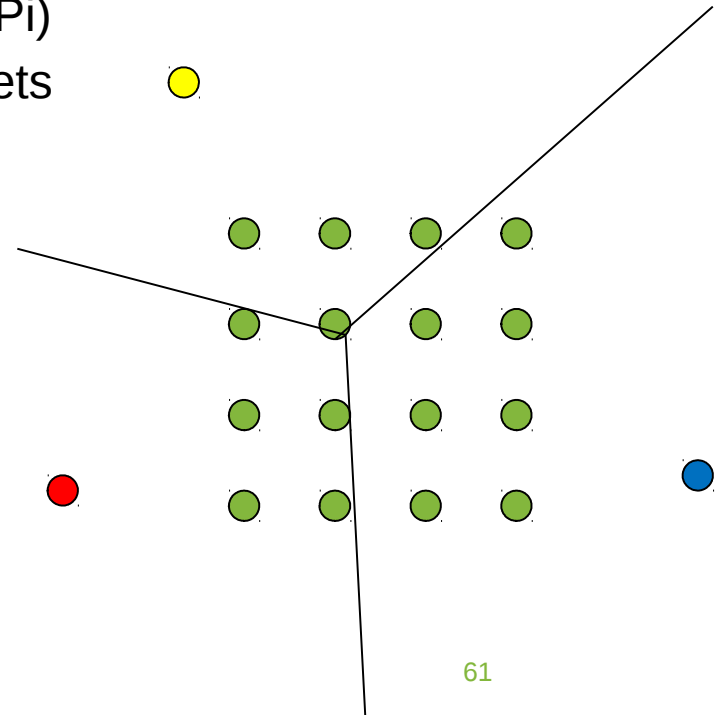
Third Technique: Reducing the Variance

A construction for $k > 1$

- We recursively partition P into subsets $\{P_i\}$ by applying the following rule:
- Let P_i be the set with the highest cost $\text{Opt}_1(P_i)$
- **If** P_i is pseudorandom **then** replace P_i by $\mu(P_i)$
- **else** find an optimal clustering of P_i into k sets $\{Q_j\}$ and replace P_i by the $\{Q_j\}$

$[\sum_j \text{Opt}_1(Q_j) \text{ decreases by a factor of } 1+\epsilon]$

- If $\sum_j \text{Opt}_1(Q_j) \leq \epsilon^2 \text{Opt}_k(P)$ then replace all $\{Q_j\}$ by $\{\mu_j\}$



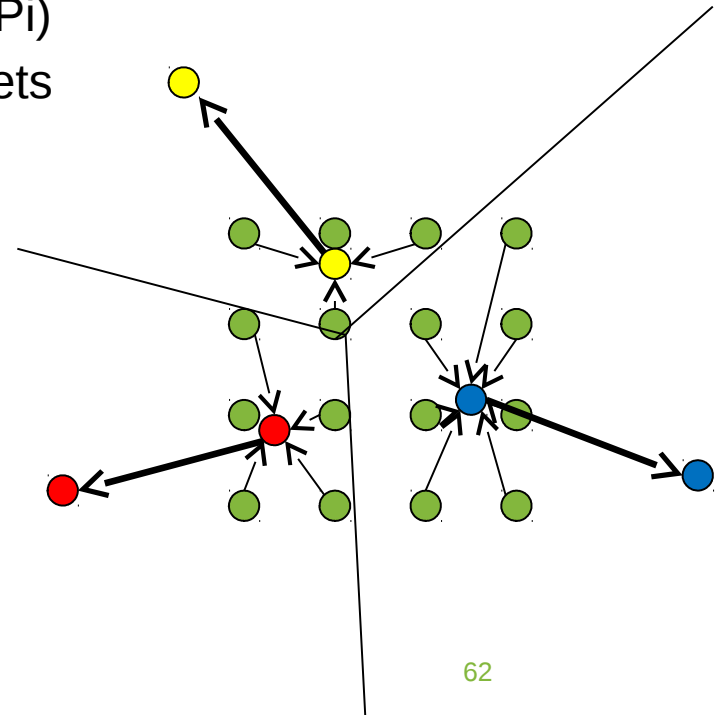
Third Technique: Reducing the Variance

A construction for $k > 1$

- We recursively partition P into subsets $\{P_i\}$ by applying the following rule:
- Let P_i be the set with the highest cost $\text{Opt}_1(P_i)$
- **If** P_i is pseudorandom **then** replace P_i by $\mu(P_i)$
- **else** find an optimal clustering of P_i into k sets $\{Q_j\}$ and replace P_i by the $\{Q_j\}$

$[\sum_j \text{Opt}_1(Q_j) \text{ decreases by a factor of } 1+\epsilon]$

- If $\sum_j \text{Opt}_1(Q_j) \leq \epsilon^2 \text{Opt}_k(P)$ then replace all $\{Q_j\}$ by $\{\mu_j\}$



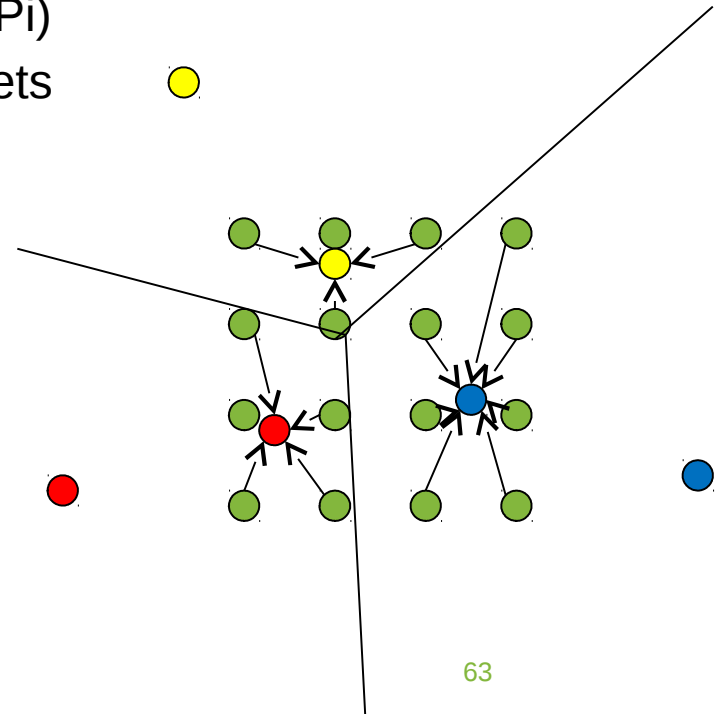
Third Technique: Reducing the Variance

A construction for $k > 1$

- We recursively partition P into subsets $\{P_i\}$ by applying the following rule:
- Let P_i be the set with the highest cost $\text{Opt}_1(P_i)$
- **If** P_i is pseudorandom **then** replace P_i by $\mu(P_i)$
- **else** find an optimal clustering of P_i into k sets $\{Q_j\}$ and replace P_i by the $\{Q_j\}$

$[\sum_j \text{Opt}_1(Q_j)$ decreases by a factor of $1+\epsilon]$

- If $\sum_j \text{Opt}_1(Q_j) \leq \epsilon^2 \text{Opt}_k(P)$ then replace all $\{Q_j\}$ by $\{\mu_j\}$



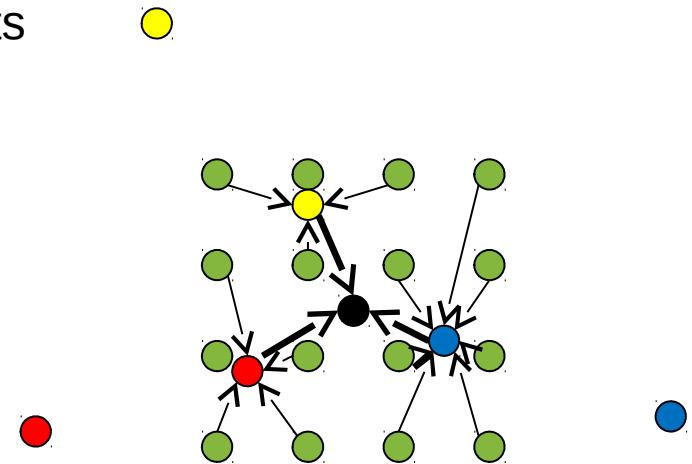
Third Technique: Reducing the Variance

A construction for $k > 1$

- We recursively partition P into subsets $\{P_i\}$ by applying the following rule:
- Let P_i be the set with the highest cost $\text{Opt}_1(P_i)$
- **If** P_i is pseudorandom **then** replace P_i by $\mu(P_i)$
- **else** find an optimal clustering of P_i into k sets $\{Q_j\}$ and replace P_i by the $\{Q_j\}$

$[\sum_j \text{Opt}_1(Q_j) \text{ decreases by a factor of } 1+\epsilon]$

- If $\sum_j \text{Opt}_1(Q_j) \leq \epsilon^2 \text{Opt}_k(P)$ then replace all $\{Q_j\}$ by $\{\mu_j\}$



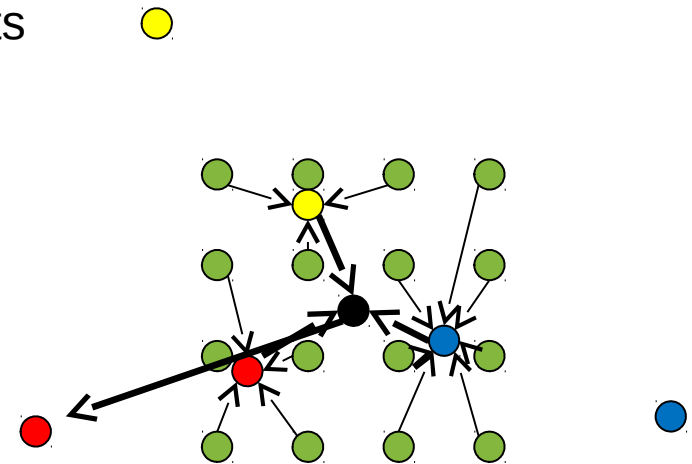
Third Technique: Reducing the Variance

A construction for $k > 1$

- We recursively partition P into subsets $\{P_i\}$ by applying the following rule:
- Let P_i be the set with the highest cost $\text{Opt}_1(P_i)$
- **If** P_i is pseudorandom **then** replace P_i by $\mu(P_i)$
- **else** find an optimal clustering of P_i into k sets $\{Q_j\}$ and replace P_i by the $\{Q_j\}$

$[\sum_j \text{Opt}_1(Q_j) \text{ decreases by a factor of } 1+\epsilon]$

- If $\sum_j \text{Opt}_1(Q_j) \leq \epsilon^2 \text{Opt}_k(P)$ then replace all $\{Q_j\}$ by $\{\mu_j\}$



Third Technique: Reducing the Variance

Theorem [**Feldman, Schmidt, S., 2013**]

- The constructed coresnet has size $k^{O(1/\epsilon \log 1/\epsilon)}$.

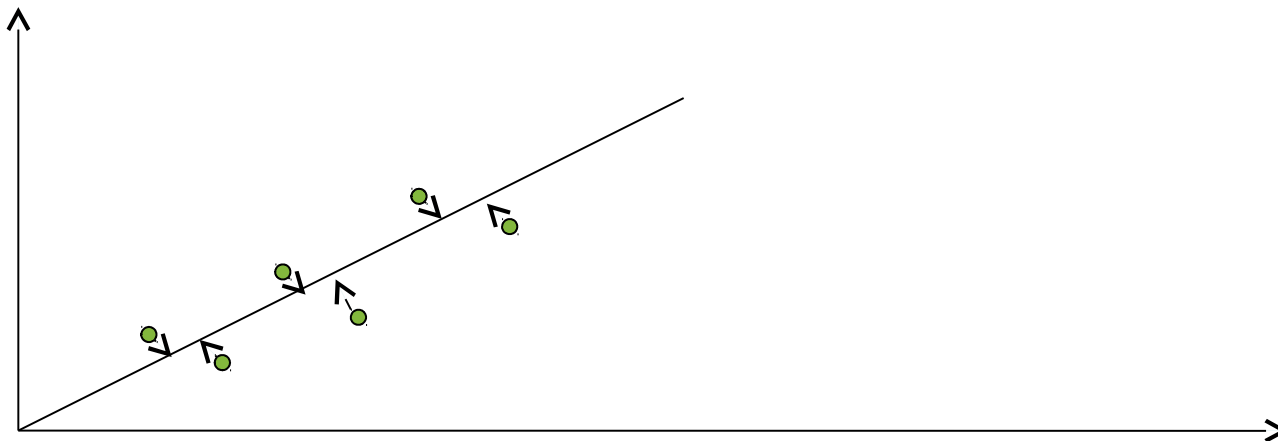
Fourth Technique: Linear Projections

Theorem [Feldman, Schmidt, S., 2013]

- There is a coresets with offset of size $(k/\epsilon)O(1)$ for the k-means problem.

Main idea

- Replace the point set by its projection on an optimal fit $O(k/\epsilon^2)$ -subspace and compute a coresets for this subspace (actually, a slightly larger space)



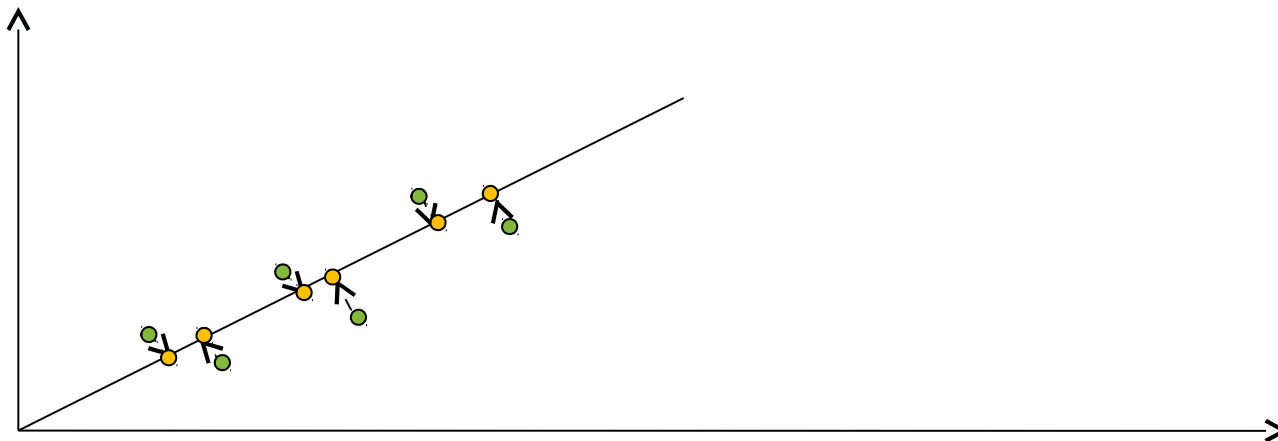
Fourth Technique: Linear Projections

Theorem [Feldman, Schmidt, S., 2013]

- There is a coresets with offset of size $(k/\epsilon)O(1)$ for the k-means problem.

Main idea

- Replace the point set by its projection on an optimal fit $O(k/\epsilon^2)$ -subspace and compute a coresets for this subspace (actually, a slightly larger space)



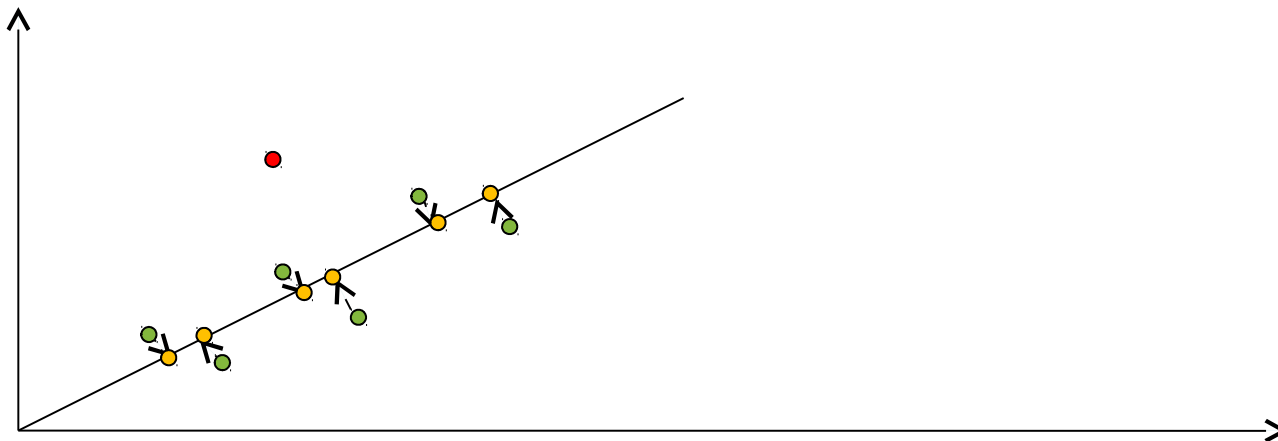
Fourth Technique: Linear Projections

Theorem [Feldman, Schmidt, S., 2013]

- There is a coresets with offset of size $(k/\epsilon)O(1)$ for the k-means problem.

Main idea

- Replace the point set by its projection on an optimal fit $O(k/\epsilon^2)$ -subspace and compute a coresets for this subspace (actually, a slightly larger space)



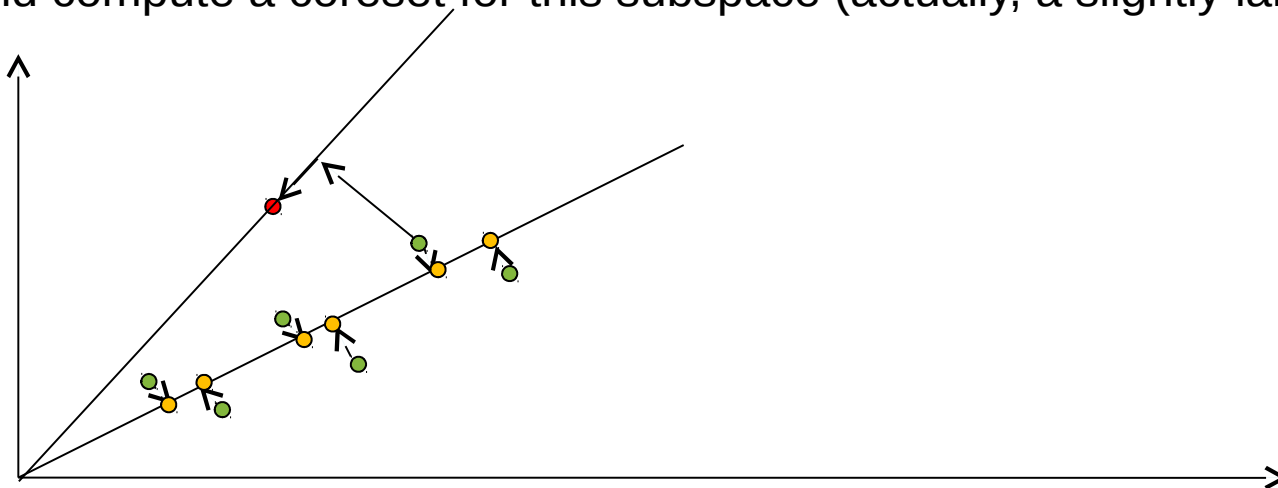
Fourth Technique: Linear Projections

Theorem [Feldman, Schmidt, S., 2013]

- There is a coresets with offset of size $(k/\epsilon)O(1)$ for the k-means problem.

Main idea

- Replace the point set by its projection on an optimal fit $O(k/\epsilon^2)$ -subspace and compute a coresets for this subspace (actually, a slightly larger space)



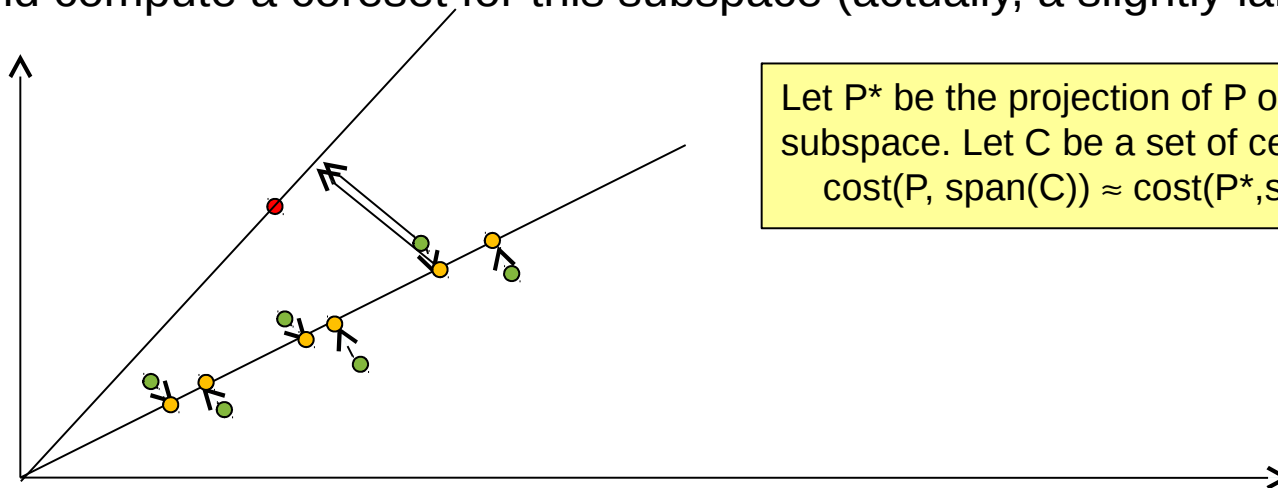
Fourth Technique: Linear Projections

Theorem [Feldman, Schmidt, S., 2013]

- There is a coresets with offset of size $(k/\epsilon)O(1)$ for the k-means problem.

Main idea

- Replace the point set by its projection on an optimal fit $O(k/\epsilon^2)$ -subspace and compute a coresets for this subspace (actually, a slightly larger space)



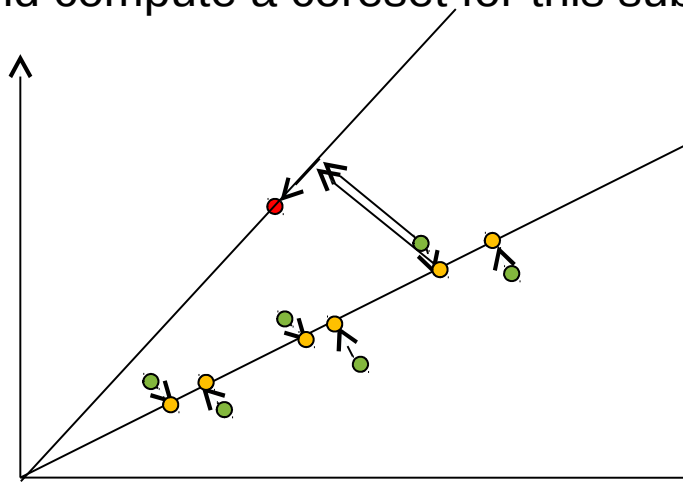
Fourth Technique: Linear Projections

Theorem [Feldman, Schmidt, S., 2013]

- There is a coresets with offset of size $(k/\epsilon)O(1)$ for the k-means problem.

Main idea

- Replace the point set by its projection on an optimal fit $O(k/\epsilon^2)$ -subspace and compute a coresets for this subspace (actually, a slightly larger space)



Let P^* be the projection of P on the best fit subspace. Let C be a set of centers. Then $\text{cost}(P, \text{span}(C)) \approx \text{cost}(P^*, \text{span}(C)) + \Delta$.

Let P^* be the projection of P on the best fit subspace. For any set of centers C , the projections of P and P^* on the linear span of C have „distance“ at most $\epsilon \text{cost}(P, \text{span}(C))$.

Extensions

Other types of centers

- Most techniques extend to different types of centers like subspace or any shapes that are contained in a low dimensional subspace

Other distance measures

- l_2^2 - distance can be replaced by l_p – distance for constant p and q

Coresets and Approximation Algorithms

Theorem

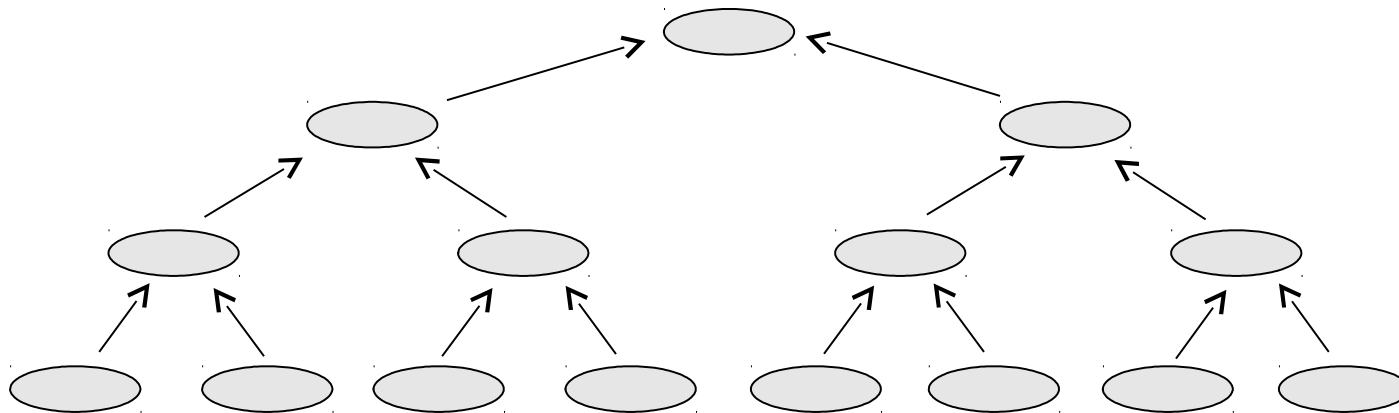
- The k -means problem can be approximated within a factor of $(1+\varepsilon)$ in time $O(T(n,k,\varepsilon) + S(n,k,\varepsilon)O(dk^2))$, where
- $T(n,k,\varepsilon)$ is the time to compute a (k,ε) -coreset on an input of size n and
- $S(n,k,\varepsilon)$ is the size of the coreset.

Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

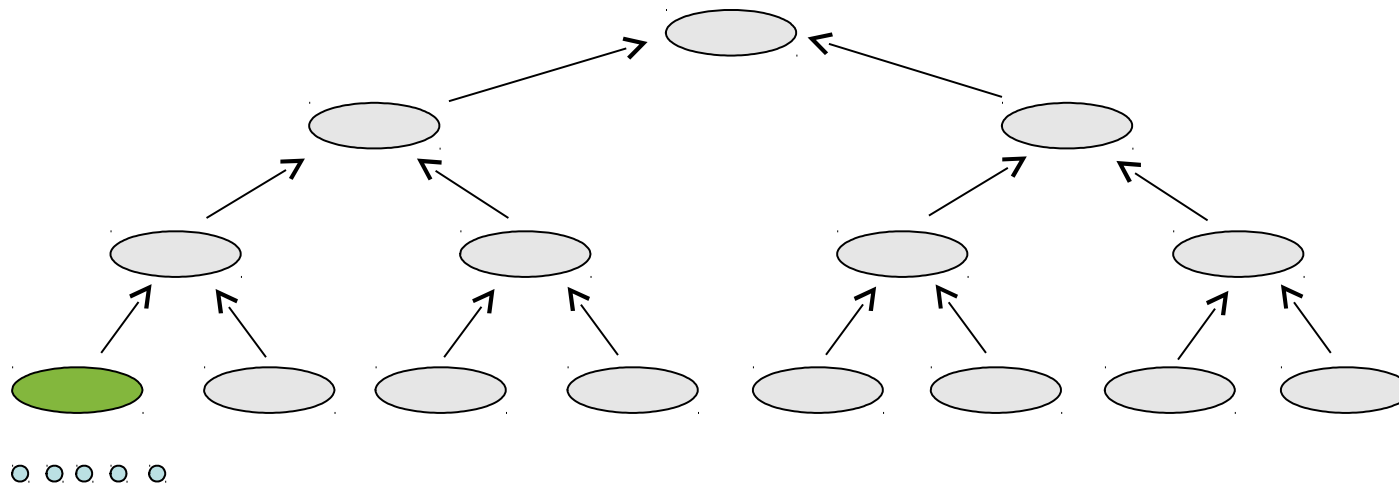


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

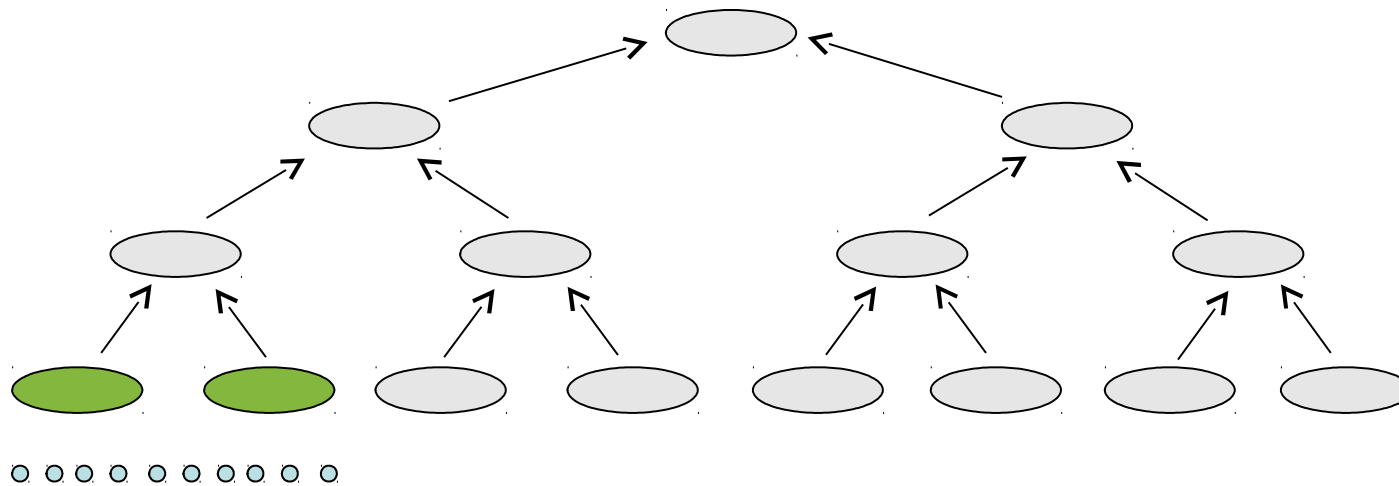


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

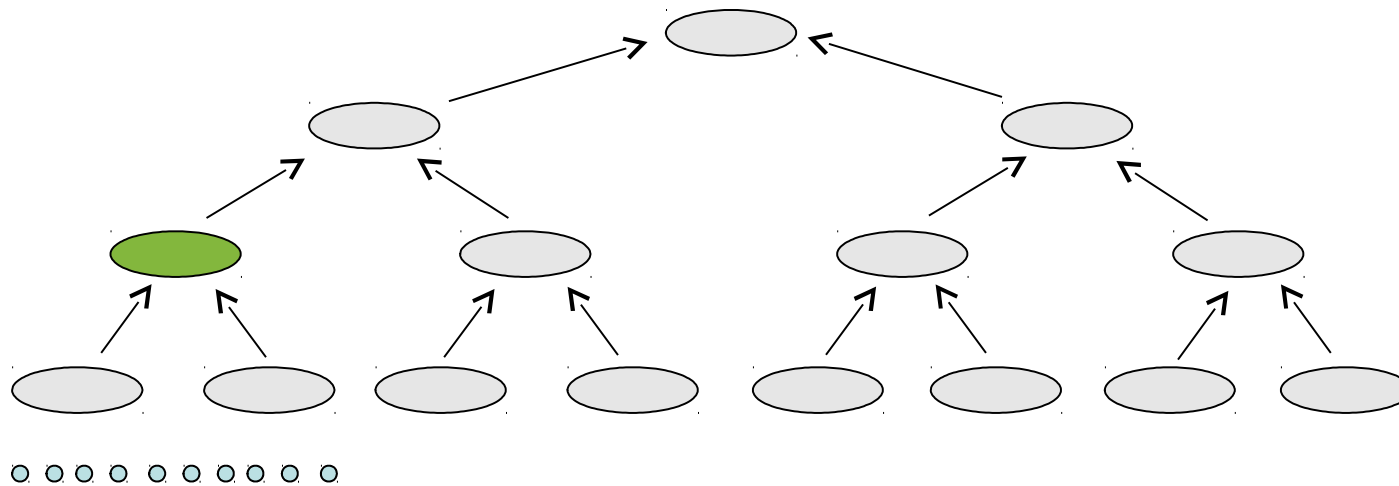


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

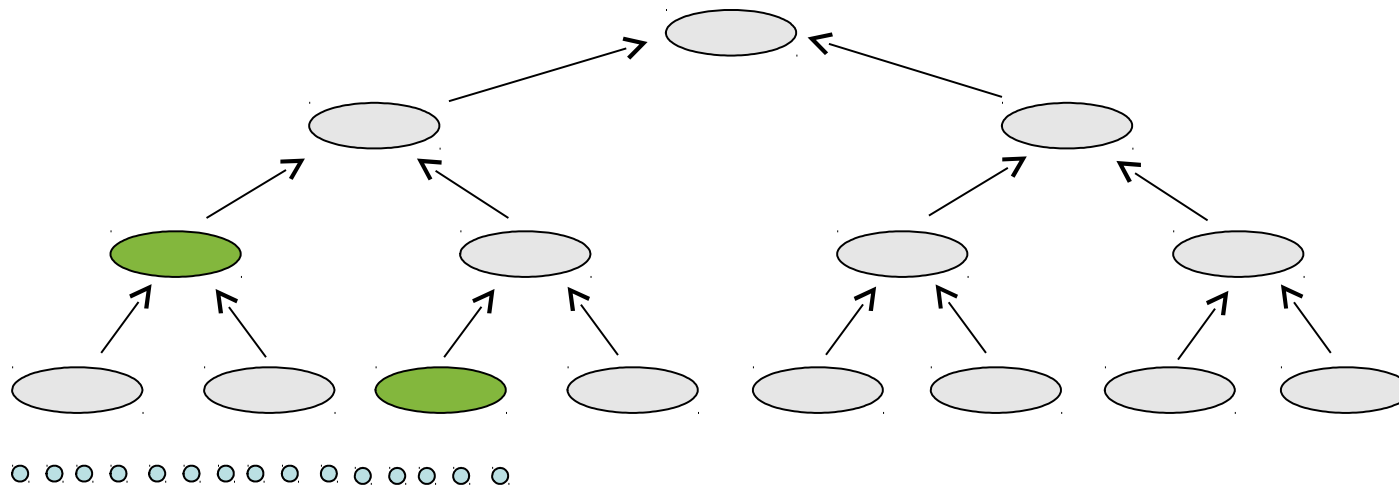


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

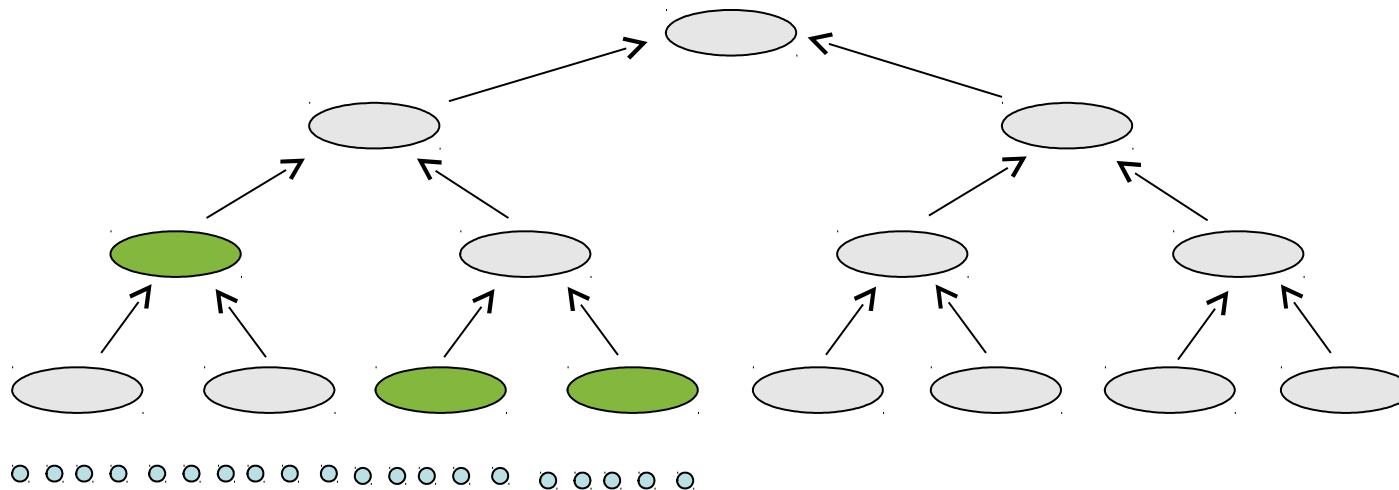


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

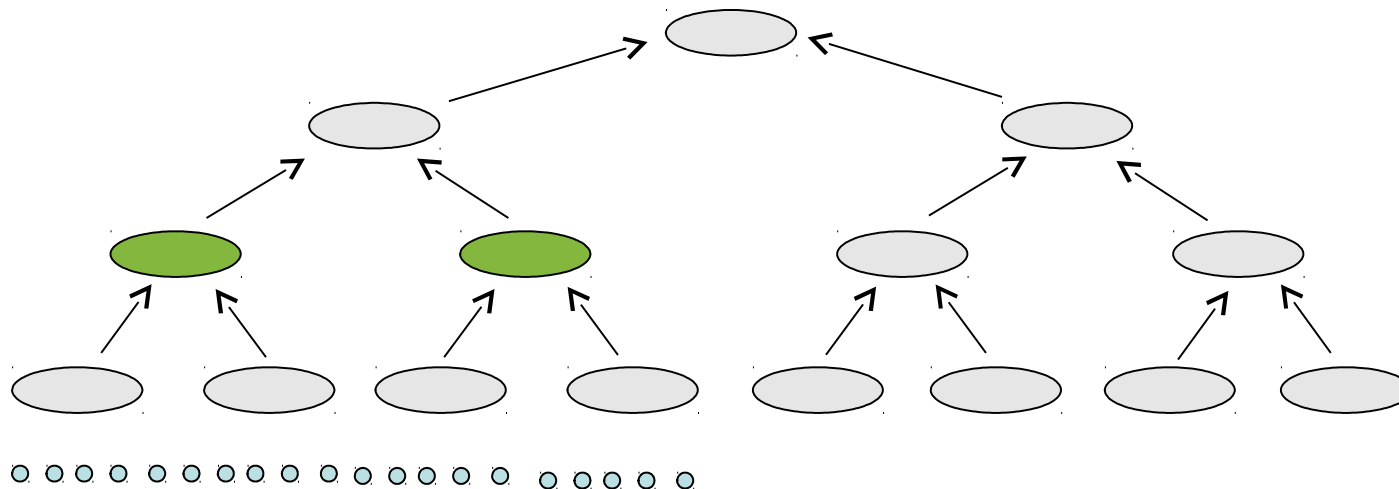


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

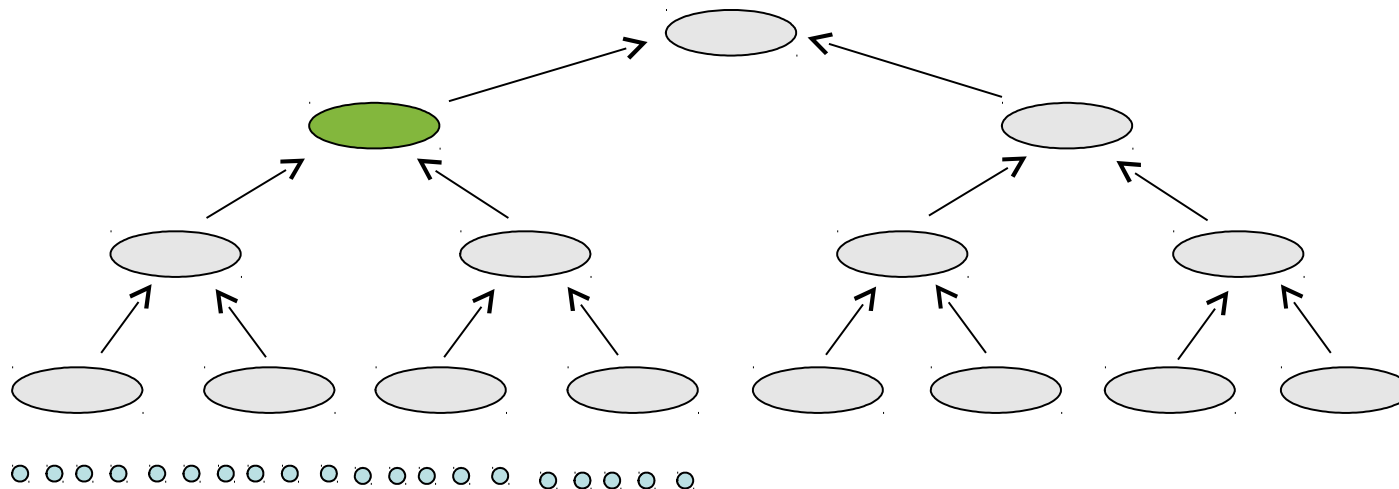


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

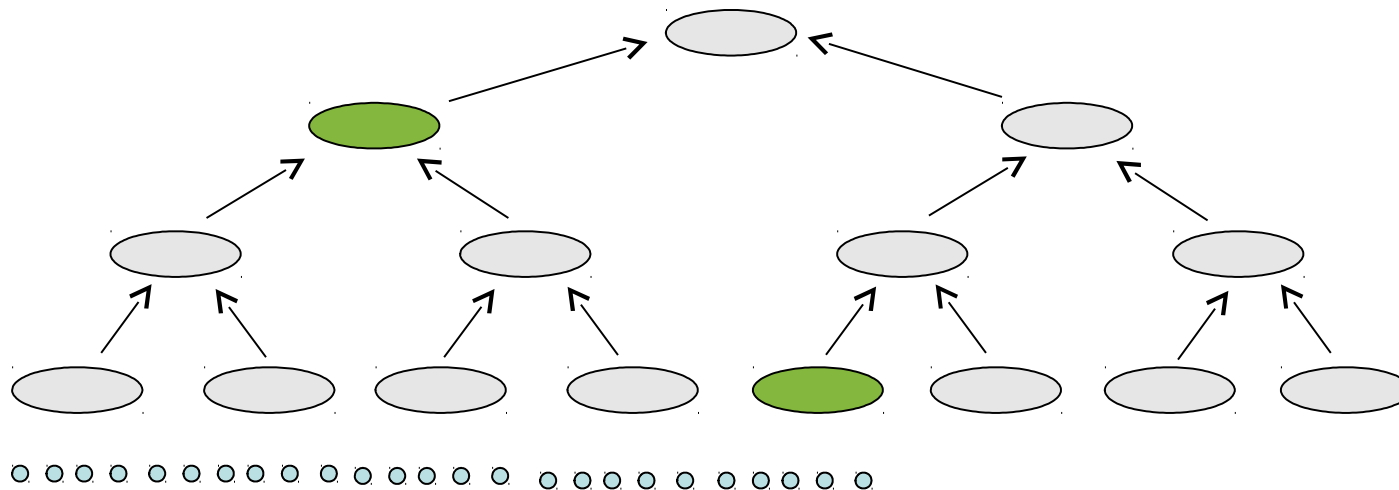


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

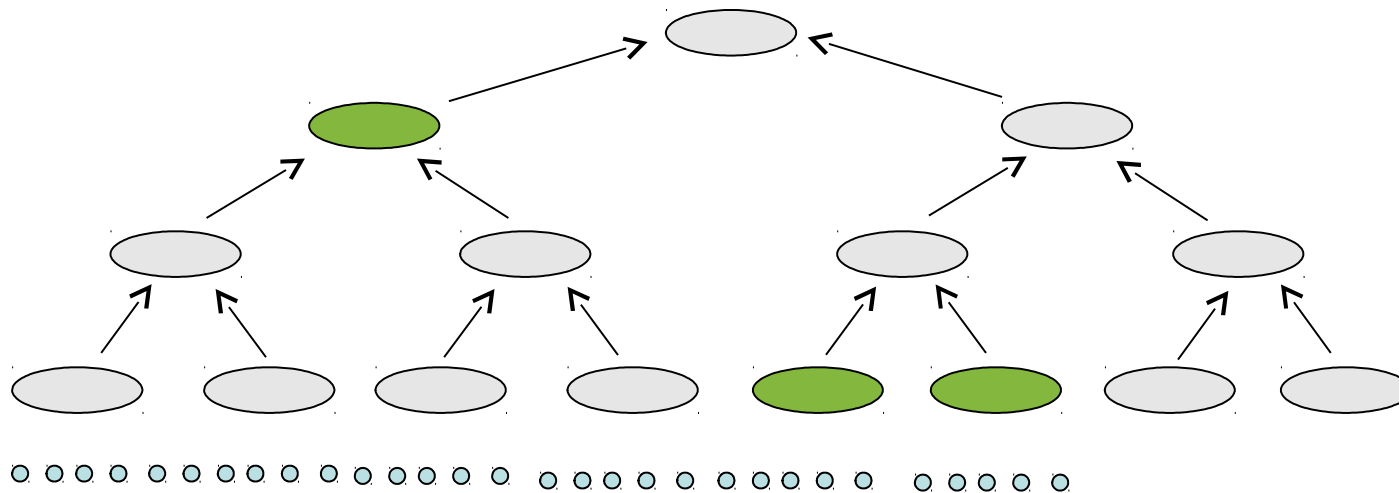


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

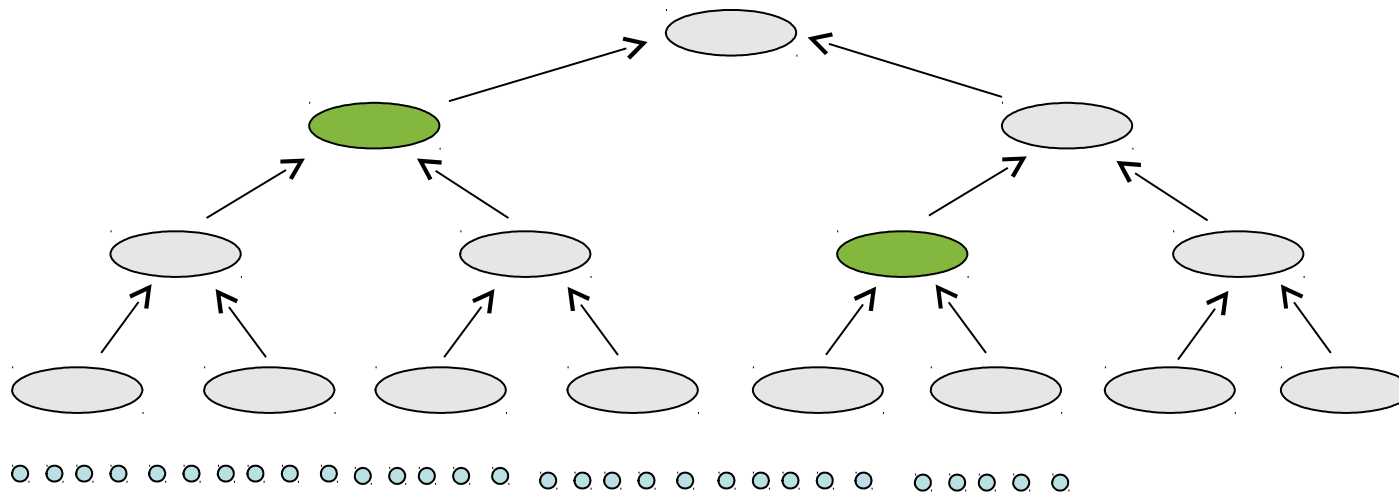


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

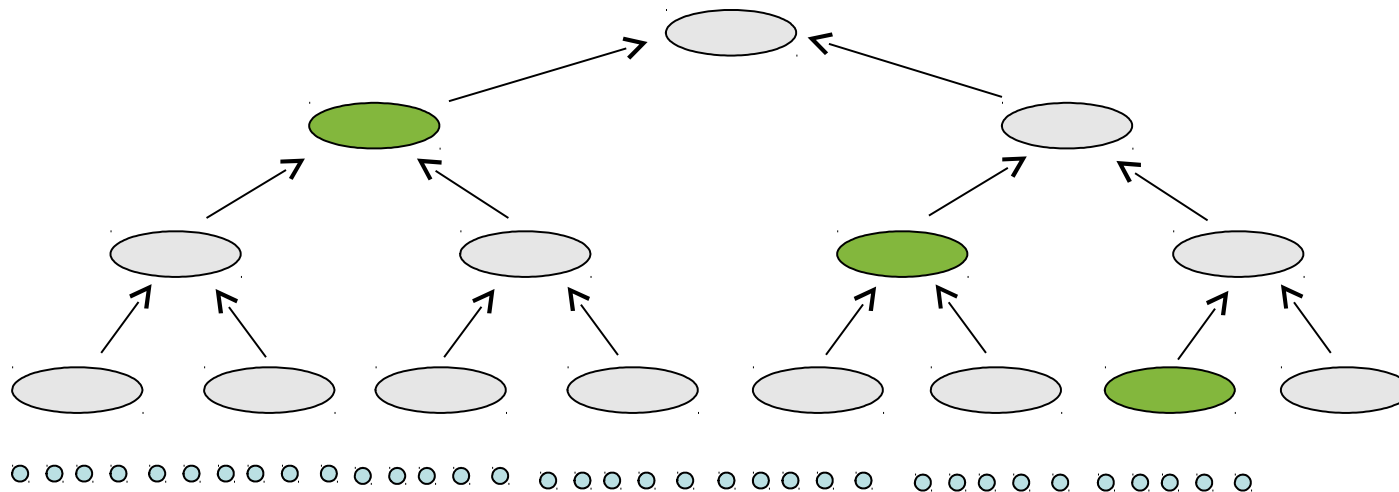


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

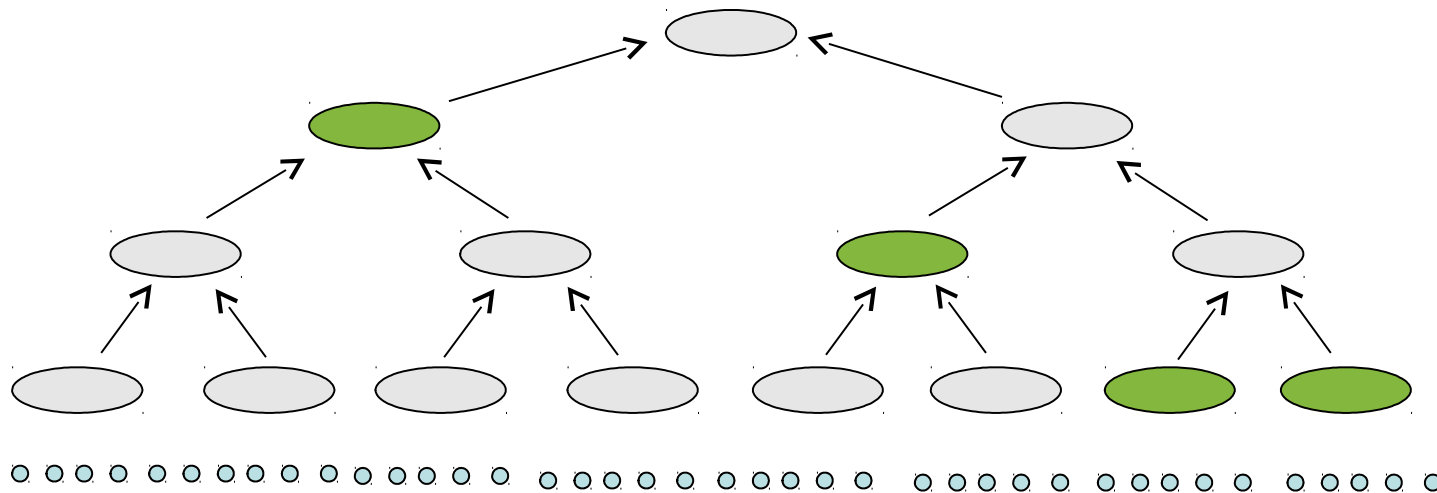


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

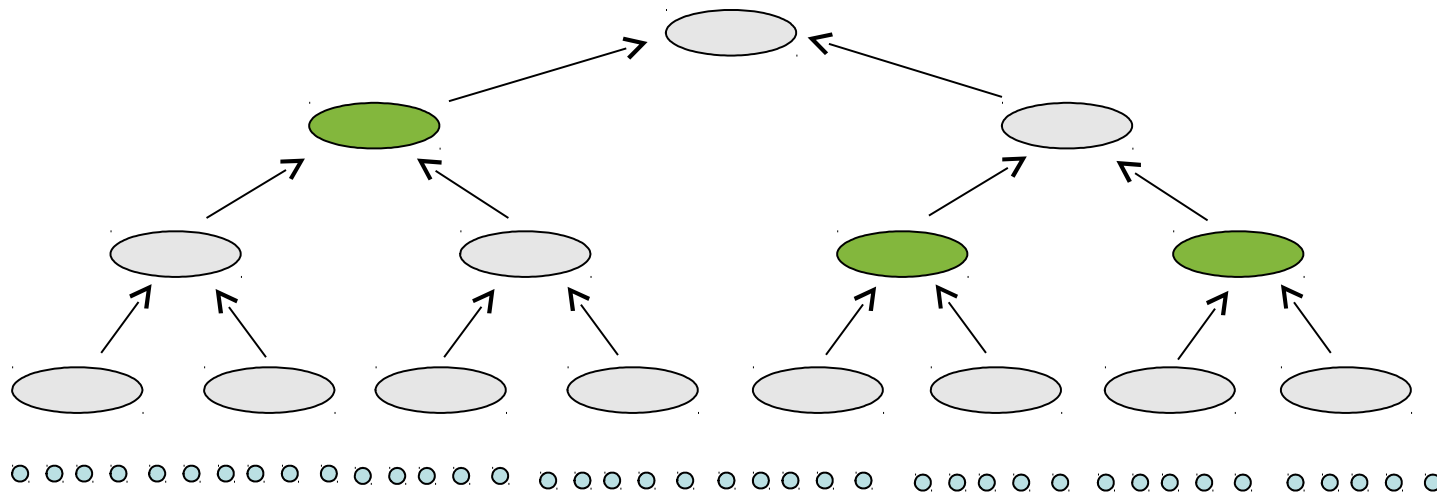


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

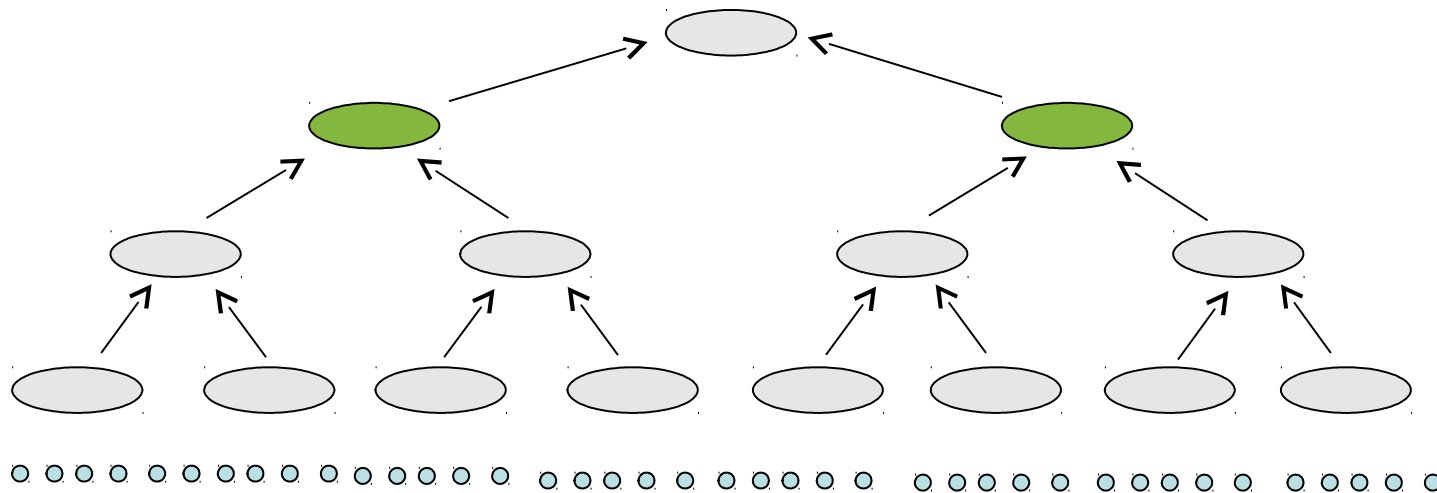


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

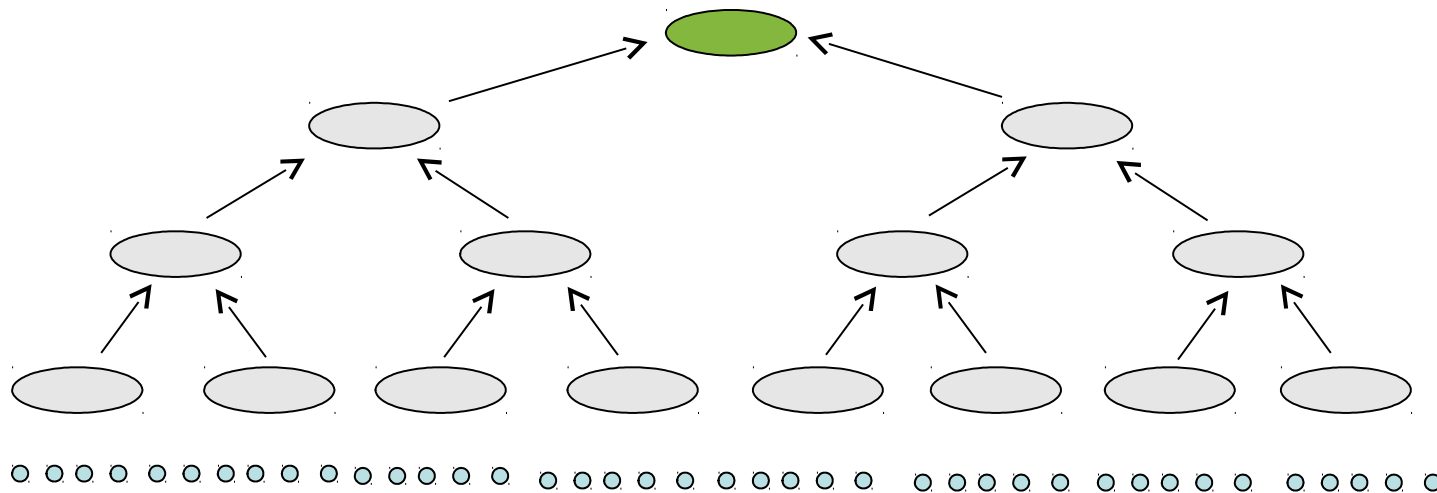


Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset



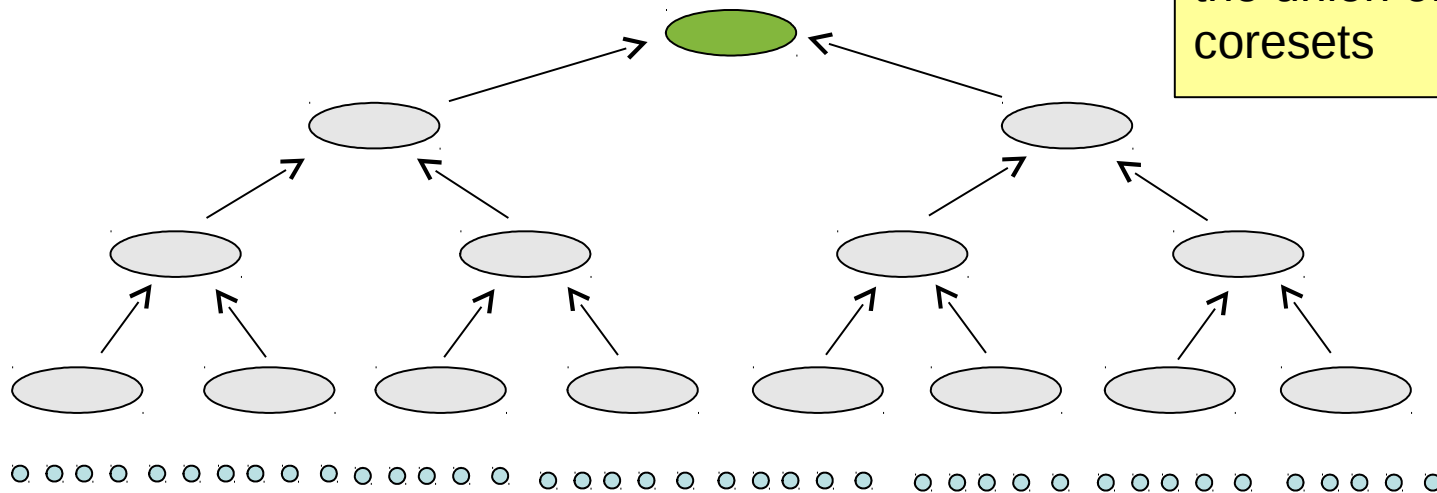
Coresets and Streaming Algorithms

[Bentley, Saxe, 80; Agarwal, Har-Peled, Varadarajan, 2004]

Assumption

- The union of two coresets for A and B is a coreset of $A \cup B$
- It is possible to compute a coreset of a coreset

Compute the clustering on the union of the stored coresets



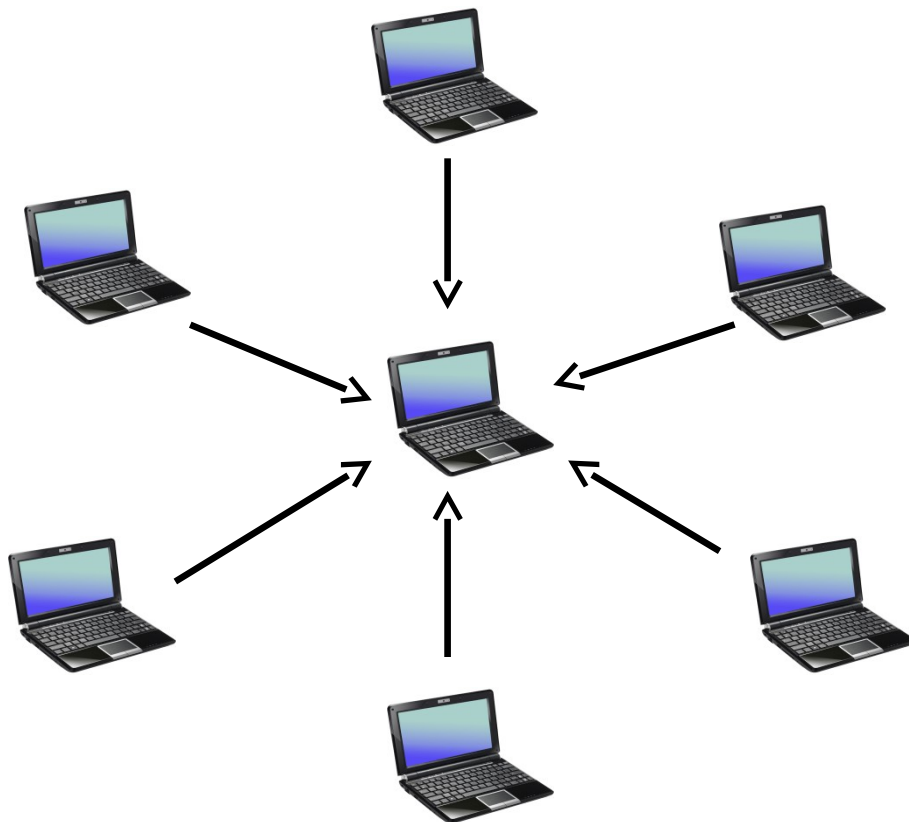
Coresets and Distributed Algorithms / Cloud Computing



Scenario: Distributed Input Data

- Data is collected locally
- Coreset will be computed locally

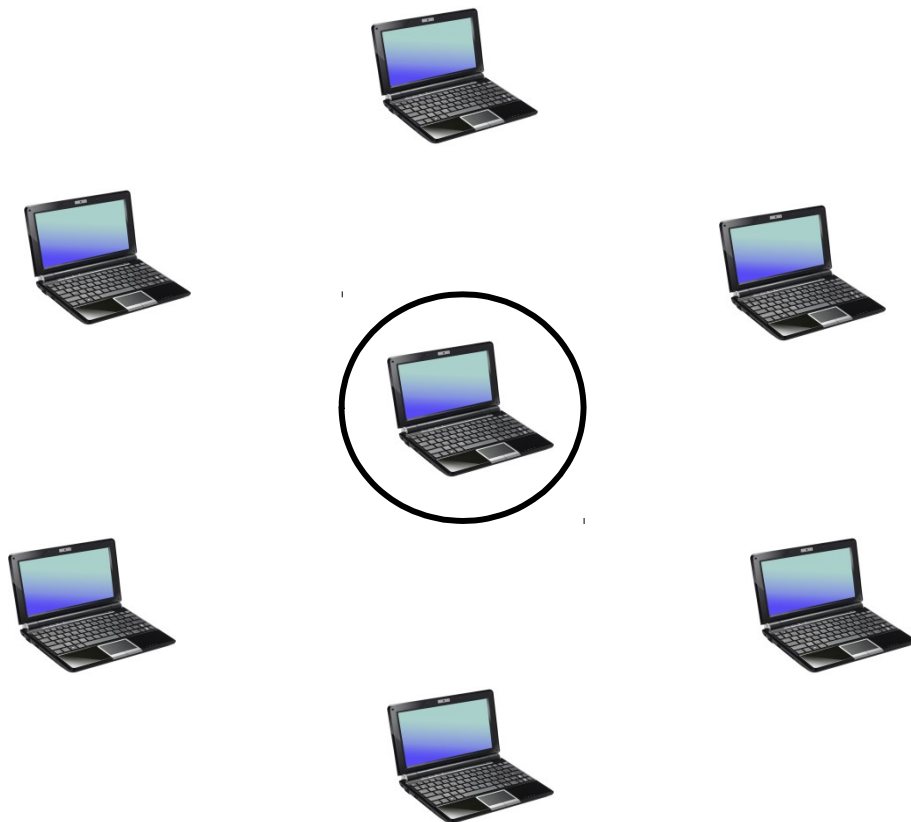
Coresets and Distributed Algorithms / Cloud Computing



Scenario: Distributed Input Data

- Data is collected locally
- Coreset will be computed locally
- Only the coresets are sent to a central computer

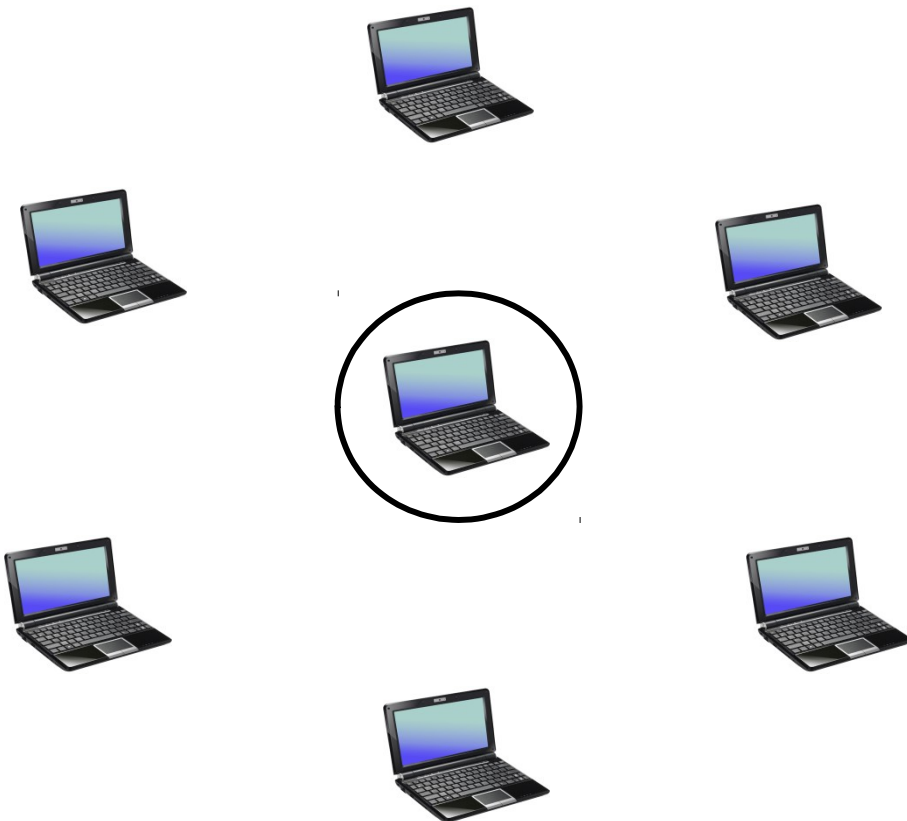
Coresets and Distributed Algorithms / Cloud Computing



Scenario: Distributed Input Data

- Data is collected locally
- Coreset will be computed locally
- Only the coresets are sent to a central computer
- Compute the clustering on this single machine

Coresets and Distributed Algorithms / Cloud Computing



Scenario: Distributed Input Data

- Data is collected locally
- Coreset will be computed locally
- Only the coresets are sent to a central computer
- Compute the clustering on this single machine

Improvements

- Implement merge&reduce like structure to avoid congestion
- Distributed computation of clustering

Practicality

StreamKM++

[Ackermann, Märtens, Raupach, Swierkot, Lammersen, S., 2011]

- Merge&Reduce based algorithm
- Coresets based on k-means++ seeding
- Coreset tree: Allows fast approximate k-means++ sampling

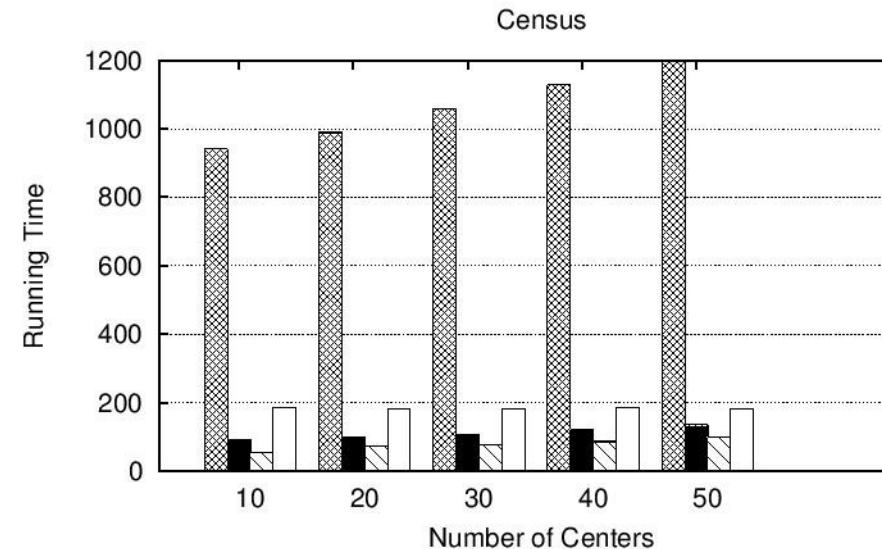
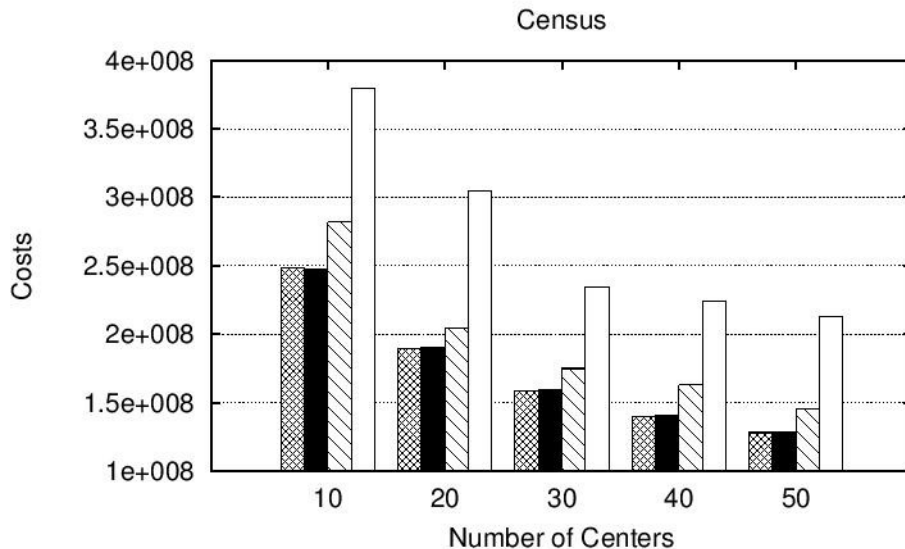
BICO [Fichtenberger, Gillé, Schmidt, Schwiegelshohn, S., 2013]

- Based on BIRCH's architecture
- Modifications, so that the algorithm computes a coreset
- A lot of heuristics to speed up performance

Practicality

Comparison of streaming clustering algorithms

- Census data set from UCI library
- 2.4 mio data points
- 68 dimensions



Summary

Coresets

- Bounded movement
- Sampling / Importance sampling
- Variance reduction
- Linear projections
- Easy-to-implement streaming & distributed algorithms
- High practicality