

Smaller Coresets for *k*-Median and *k*-Means Clustering

Sariel Har-Peled Akash Kushal
UIUC, Urbana, IL

Coresets for Clustering

Clustering

P: set of **n** points in \mathbb{R}^d , compute a set **C** of **k** centers –

- **k-median**: sum dist. points in **P** to nearest center.

$$v_{\mathbf{C}}(\mathbf{P}) = \sum_{p \in \mathbf{P}} d(p, \mathbf{C})$$

- **k-means**: sum squares dist. points in **P** to nearest center.

Target: Find **C** that minimizes price.

Coreset

Definition (Coreset)

$\mathcal{S} \subseteq \mathbf{P}$ - A weighted subset is (k, ϵ) -coreset if

$$\forall \mathbf{C} \subseteq \mathbb{R}^d \quad \nu_{\mathbf{C}}(\mathbf{P}) \approx_{\epsilon} \nu_{\mathbf{C}}(\mathcal{S}),$$

where $\mathbf{a} \approx_{\epsilon} \mathbf{b}$ if $(1 - \epsilon)\mathbf{a} \leq \mathbf{b} \leq (1 + \epsilon)\mathbf{a}$.

Coreset = small sketch of the input for clustering.

Question

What is the smallest sketch possible?

Coreset

Definition (Coreset)

$\mathcal{S} \subseteq \mathbf{P}$ - A weighted subset is (k, ϵ) -coreset if

$$\forall \mathbf{C} \subseteq \mathbb{R}^d \quad \nu_{\mathbf{C}}(\mathbf{P}) \approx_{\epsilon} \nu_{\mathbf{C}}(\mathcal{S}),$$

where $\mathbf{a} \approx_{\epsilon} \mathbf{b}$ if $(1 - \epsilon)\mathbf{a} \leq \mathbf{b} \leq (1 + \epsilon)\mathbf{a}$.

Coreset = small sketch of the input for clustering.

Question

What is the smallest sketch possible?

Coreset

Definition (Coreset)

$\mathcal{S} \subseteq \mathbf{P}$ - A weighted subset is (k, ϵ) -coreset if

$$\forall \mathbf{C} \subseteq \mathbb{R}^d \quad \nu_{\mathbf{C}}(\mathbf{P}) \approx_{\epsilon} \nu_{\mathbf{C}}(\mathcal{S}),$$

where $\mathbf{a} \approx_{\epsilon} \mathbf{b}$ if $(1 - \epsilon)\mathbf{a} \leq \mathbf{b} \leq (1 + \epsilon)\mathbf{a}$.

Coreset = small sketch of the input for clustering.

Question

What is the smallest sketch possible?

k-median

- **[Arora et al., 1998]:**
 $O(n^{O(1/\epsilon)})$ time $(1 + \epsilon)$ -approximation for points in plane
- **[Kolliopoulos and Rao, 1999]** (Discrete version):
 $O(\varrho n \log n \log k)$ time, where $\varrho = \exp\left[O(\epsilon^{-1} \log 1/\epsilon)^{d-1}\right]$.
- **[Har-Peled and Mazumdar, 2004]:**
 Coreset of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 Improve running time to $O\left(n + \varrho k^{O(1)} \log^{O(1)} n\right)$

New Result:

Show coreset of size $O(k^2/\epsilon^d)$

Independent of $n!$

k-median

- **[Arora et al., 1998]:**
 $O(n^{O(1/\epsilon)})$ time $(1 + \epsilon)$ -approximation for points in plane
- **[Kolliopoulos and Rao, 1999]** (Discrete version):
 $O(\varrho n \log n \log k)$ time, where $\varrho = \exp\left[O(\epsilon^{-1} \log 1/\epsilon)^{d-1}\right]$.
- **[Har-Peled and Mazumdar, 2004]:**
 Coreset of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 Improve running time to $O\left(n + \varrho k^{O(1)} \log^{O(1)} n\right)$

New Result:

Show coreset of size $O(k^2/\epsilon^d)$

Independent of $n!$

k-median

- **[Arora et al., 1998]:**
 $O(n^{O(1/\epsilon)})$ time $(1 + \epsilon)$ -approximation for points in plane
- **[Kolliopoulos and Rao, 1999]** (Discrete version):
 $O(\varrho n \log n \log k)$ time, where $\varrho = \exp\left[O(\epsilon^{-1} \log 1/\epsilon)^{d-1}\right]$.
- **[Har-Peled and Mazumdar, 2004]:**
 Coreset of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 Improve running time to $O\left(n + \varrho k^{O(1)} \log^{O(1)} n\right)$

New Result:

Show coreset of size $O(k^2/\epsilon^d)$

Independent of $n!$

k-median

- **[Arora et al., 1998]:**
 $O(n^{O(1/\epsilon)})$ time $(1 + \epsilon)$ -approximation for points in plane
- **[Kolliopoulos and Rao, 1999]** (Discrete version):
 $O(\varrho n \log n \log k)$ time, where $\varrho = \exp\left[O(\epsilon^{-1} \log 1/\epsilon)^{d-1}\right]$.
- **[Har-Peled and Mazumdar, 2004]:**
 Coreset of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 Improve running time to $O\left(n + \varrho k^{O(1)} \log^{O(1)} n\right)$

New Result:

Show coreset of size $O(k^2/\epsilon^d)$

Independent of **n!**

New coresets for *k*-median of size

- independent of ***n***!
- disconnected from ***n***!!
- free of ***n***!!!
- autonomous of ***n***!!!!
- unregimented by ***n***!!!!!
- unblemished by ***n***!!!!!!!
- unmolested by ***n***!!!!!!!!!
- !!!

k -means

- **[de la Vega et al., 2003]:** $(1 + \epsilon)$ -approx in high dim.
Running time $O(g(k, \epsilon)dn \log^k n)$, where
 $g(k, \epsilon) = \exp[(k^3/\epsilon^8)(\ln(k/\epsilon)) \ln k]$
- **[Kumar et al., 2004]:**
improved to $O(h(k, \epsilon)dn)$, where $h(k, \epsilon) = 2^{(k/\epsilon)^{O(1)}}$
- **[Matoušek, 2000]:** running time $O_{d,k}(n\epsilon^{-2k^2d} \log^k n)$
- **[Har-Peled and Mazumdar, 2004]:**
Coresets of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 ϵ -approx in $O(n + k^{k+2}\epsilon^{-(2d+1)k} \log^{k+1}(n/\epsilon))$ time
- **New Result:** show coreset of size $O\left(k^3/\epsilon^{d+1}\right)$

k-means

- **[de la Vega et al., 2003]:** $(1 + \epsilon)$ -approx in high dim.
Running time $O(g(k, \epsilon)dn \log^k n)$, where
 $g(k, \epsilon) = \exp[(k^3/\epsilon^8)(\ln(k/\epsilon)) \ln k]$
- **[Kumar et al., 2004]:**
improved to $O(h(k, \epsilon)dn)$, where $h(k, \epsilon) = 2^{(k/\epsilon)^{O(1)}}$
- **[Matoušek, 2000]:** running time $O_{d,k}(n\epsilon^{-2k^2d} \log^k n)$
- **[Har-Peled and Mazumdar, 2004]:**
Coresets of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 ϵ -approx in $O(n + k^{k+2}\epsilon^{-(2d+1)k} \log^{k+1}(n/\epsilon))$ time
- **New Result:** show coreset of size $O\left(k^3/\epsilon^{d+1}\right)$

k -means

- **[de la Vega et al., 2003]**: $(1 + \epsilon)$ -approx in high dim.
Running time $O(g(k, \epsilon)dn \log^k n)$, where
 $g(k, \epsilon) = \exp[(k^3/\epsilon^8)(\ln(k/\epsilon)) \ln k]$
- **[Kumar et al., 2004]**:
improved to $O(h(k, \epsilon)dn)$, where $h(k, \epsilon) = 2^{(k/\epsilon)^{O(1)}}$
- **[Matoušek, 2000]**: running time $O_{d,k}(n\epsilon^{-2k^2d} \log^k n)$
- **[Har-Peled and Mazumdar, 2004]**:
Coresets of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 ϵ -approx in $O(n + k^{k+2} \epsilon^{-(2d+1)k} \log^{k+1}(n/\epsilon))$ time
- **New Result**: show coreset of size $O\left(k^3/\epsilon^{d+1}\right)$

k -means

- **[de la Vega et al., 2003]**: $(1 + \epsilon)$ -approx in high dim.
Running time $O(g(k, \epsilon)dn \log^k n)$, where
 $g(k, \epsilon) = \exp[(k^3/\epsilon^8)(\ln(k/\epsilon)) \ln k]$
- **[Kumar et al., 2004]**:
improved to $O(h(k, \epsilon)dn)$, where $h(k, \epsilon) = 2^{(k/\epsilon)^{O(1)}}$
- **[Matoušek, 2000]**: running time $O_{d,k}(n\epsilon^{-2k^2d} \log^k n)$
- **[Har-Peled and Mazumdar, 2004]**:
Coresets of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 ϵ -approx in $O(n + k^{k+2}\epsilon^{-(2d+1)k} \log^{k+1}(n/\epsilon))$ time
- **New Result**: show coreset of size $O(k^3/\epsilon^{d+1})$

k -means

- **[de la Vega et al., 2003]:** $(1 + \epsilon)$ -approx in high dim.
Running time $O(g(k, \epsilon)dn \log^k n)$, where
 $g(k, \epsilon) = \exp[(k^3/\epsilon^8)(\ln(k/\epsilon)) \ln k]$
- **[Kumar et al., 2004]:**
improved to $O(h(k, \epsilon)dn)$, where $h(k, \epsilon) = 2^{(k/\epsilon)^{O(1)}}$
- **[Matoušek, 2000]:** running time $O_{d,k}(n\epsilon^{-2k^2d} \log^k n)$
- **[Har-Peled and Mazumdar, 2004]:**
Coresets of size $O\left(\frac{k \log n}{\epsilon^d}\right)$
 ϵ -approx in $O(n + k^{k+2}\epsilon^{-(2d+1)k} \log^{k+1}(n/\epsilon))$ time
- **New Result:** show coreset of size $O(k^3/\epsilon^{d+1})$

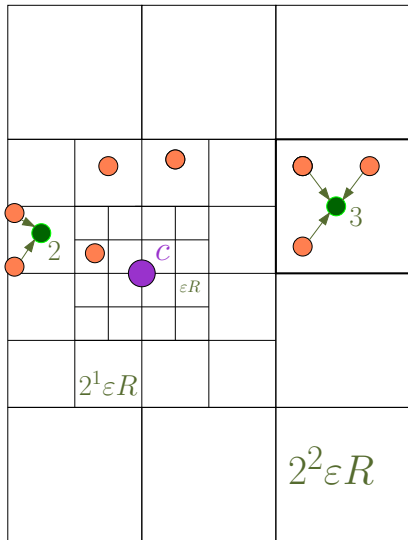
k -means - Effros & Schulman result

- Deterministic clustering with data nets
[Effros and Schulman, 2003]
- For k means only.
- Show the existence of centroid set of size independent of n .
- Complicated.
- Inspiration to conjecture that small coresets exists.
- New result implies their result (better parameters).
- New result much simpler.
- Effros & Schulman recently showed existence of small coresets for k -means.

Old Coreset Construction

Construction

- D - a set of k centers
 $\nu_D(\mathbf{P}) = O(\nu_{\text{opt}})$.
- build exponential grid of $O(\log n)$ levels around each center,
- snap the points of \mathbf{P} to this grid.
- *price of snapping is smaller than $\epsilon \nu_{\text{opt}}(\mathbf{P}, k)$*



Old Coreset Construction – limitations

Lemma \Rightarrow Problem

If snapping is “cheap” ($\leq \epsilon \nu_{\text{opt}}(\mathbf{P}, k)$), then $|\mathcal{S}| = \Omega(\log n)$.

Conclusions

For \mathcal{S} to be small

\Rightarrow Must pick \mathcal{S} such that snapping errors *cancel each other out*.

Old Coreset Construction – limitations

Lemma \Rightarrow Problem

If snapping is “cheap” ($\leq \epsilon \nu_{\text{opt}}(\mathbf{P}, k)$), then $|\mathcal{S}| = \Omega(\log n)$.

Conclusions

For \mathcal{S} to be small

\Rightarrow Must pick \mathcal{S} such that snapping errors *cancel each other out*.

Warmup exercise - Coreset in One Dimension

Problem

$P \subseteq \mathbb{R}$ - n points on real line.

Pick (k, ϵ) -coreset for P for centers on the line.



Basic Idea:

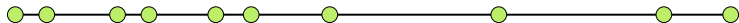
Break the point set into batches, and use the mean point of batch, as the representative for the coreset.

Warmup exercise - Coreset in One Dimension

Problem

$P \subseteq \mathbb{R}$ - n points on real line.

Pick (k, ϵ) -coreset for P for centers on the line.



Basic Idea:

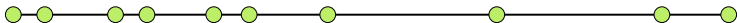
Break the point set into batches, and use the mean point of batch, as the representative for the coreset.

Warmup exercise - Coreset in One Dimension

Problem

$P \subseteq \mathbb{R}$ - n points on real line.

Pick (k, ϵ) -coreset for P for centers on the line.



Basic Idea:

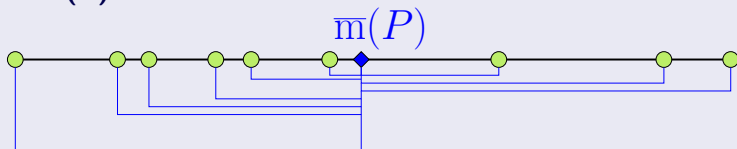
Break the point set into batches, and use the mean point of batch, as the representative for the coreset.

Warmup exercise - Coreset in One Dimension - cont'd

Definition (Cumulative Error)

$$\mathcal{E}_\nu(\mathbf{P}) = \nu_{\bar{\mathbf{m}}(\mathbf{P})}(\mathbf{P}) = \sum_{\mathbf{p} \in \mathbf{P}} w_{\mathbf{p}} \|\mathbf{p} - \bar{\mathbf{m}}\|$$

$\bar{\mathbf{m}} = \bar{\mathbf{m}}(\mathbf{P})$ is mean of \mathbf{P} .



$\mathcal{E}_\nu(\mathbf{P}) =$ Total length of blue segments.

Warmup exercise - Coreset in One Dimension - cont'd

- \mathbf{V} : $\nu_{\text{opt}}(\mathbf{P}, k) \leq \mathbf{V} \leq c \cdot \nu_{\text{opt}}(\mathbf{P}, k)$
- Sweep \mathbf{P} from left to right.
- Break into batches \mathbf{B}_i with $\mathcal{E}_\nu(\mathbf{B}_i) = \frac{\varepsilon}{10ck} \mathbf{V}$
(allow fractional points)
- # batches is $\alpha = \mathbf{O}(k/\varepsilon)$
- Coreset point: $(\bar{\mathbf{m}}(\mathbf{B}_i), \mathbf{w}(\mathbf{B}_i))$
- Coreset: $\left\{ (\bar{\mathbf{m}}(\mathbf{B}_i), \mathbf{w}(\mathbf{B}_i)) \mid i = 1, \dots, \alpha \right\}$.

Correctness

- \mathbf{B} - a batch of points
- $\mathcal{I}(\mathbf{B})$ - interval containing \mathbf{B} .
- $\forall \mathbf{c}$ outside interval $\mathcal{I}(\mathbf{B})$ then $\nu_{\mathbf{c}}(\mathbf{P}) = \nu_{\mathbf{c}}(\mathcal{S})$.
- *problematic* batches = contribute to the error
 - Batches served by two or more centers ($\leq k - 1$ batches)
 - batches with center point inside them ($\leq k$ such batches).
- Error of problematic batches $\leq \mathbf{L} = \mathcal{E}_{\nu}(\mathbf{B}_i) = \frac{\epsilon}{10ck} \mathbf{V}$
- Total error $\leq (2k - 1) \cdot \mathbf{L} \leq \epsilon \nu_{\text{opt}}(\mathbf{P}, k)$.

Construction can be extended to:

- Points on line $\ell \subseteq \mathbb{R}^d$.
- Centers somewhere in \mathbb{R}^d
- Extension not trivial (but doable - details in paper)

Correctness

- \mathbf{B} - a batch of points
- $\mathcal{I}(\mathbf{B})$ - interval containing \mathbf{B} .
- $\forall \mathbf{c}$ outside interval $\mathcal{I}(\mathbf{B})$ then $\nu_{\mathbf{c}}(\mathbf{P}) = \nu_{\mathbf{c}}(\mathcal{S})$.
- *problematic* batches = contribute to the error
 - Batches served by two or more centers ($\leq k - 1$ batches)
 - batches with center point inside them ($\leq k$ such batches).
- Error of problematic batches $\leq L = \mathcal{E}_{\nu}(\mathbf{B}_i) = \frac{\epsilon}{10ck} \mathbf{V}$
- Total error $\leq (2k - 1) \cdot L \leq \epsilon \nu_{\text{opt}}(\mathbf{P}, k)$.

Construction can be extended to:

- Points on line $\ell \subseteq \mathbb{R}^d$.
- Centers somewhere in \mathbb{R}^d
- Extension not trivial (but doable - details in paper)

Correctness

- \mathbf{B} - a batch of points
- $\mathcal{I}(\mathbf{B})$ - interval containing \mathbf{B} .
- $\forall \mathbf{c}$ outside interval $\mathcal{I}(\mathbf{B})$ then $\nu_{\mathbf{c}}(\mathbf{P}) = \nu_{\mathbf{c}}(\mathcal{S})$.
- *problematic* batches = contribute to the error
 - Batches served by two or more centers ($\leq k - 1$ batches)
 - batches with center point inside them ($\leq k$ such batches).
- Error of problematic batches $\leq \mathbf{L} = \mathcal{E}_{\nu}(\mathbf{B}_i) = \frac{\epsilon}{10ck} \mathbf{V}$
- Total error $\leq (2k - 1) \cdot \mathbf{L} \leq \epsilon \nu_{\text{opt}}(\mathbf{P}, k)$.

Construction can be extended to:

- Points on line $\ell \subseteq \mathbb{R}^d$.
- Centers somewhere in \mathbb{R}^d
- Extension not trivial (but doable - details in paper)

Correctness

- \mathbf{B} - a batch of points
- $\mathcal{I}(\mathbf{B})$ - interval containing \mathbf{B} .
- $\forall \mathbf{c}$ outside interval $\mathcal{I}(\mathbf{B})$ then $\nu_{\mathbf{c}}(\mathbf{P}) = \nu_{\mathbf{c}}(\mathcal{S})$.
- *problematic* batches = contribute to the error
 - Batches served by two or more centers ($\leq k - 1$ batches)
 - batches with center point inside them ($\leq k$ such batches).
- Error of problematic batches $\leq \mathbf{L} = \mathcal{E}_{\nu}(\mathbf{B}_i) = \frac{\epsilon}{10ck} \mathbf{V}$
- Total error $\leq (2k - 1) \cdot \mathbf{L} \leq \epsilon \nu_{\text{opt}}(\mathbf{P}, k)$.

Construction can be extended to:

- Points on line $\ell \subseteq \mathbb{R}^d$.
- Centers somewhere in \mathbb{R}^d
- Extension not trivial (but doable - details in paper)

What we currently have...

Lemma

- \mathbf{P} a set of points on a line ℓ
- ℓ - line in \mathbb{R}^d
- Coreset of size $\mathbf{O}(k/\epsilon)$ for \mathbf{P}
- For \mathbf{k} -median clustering.

Problem

General point-set do not lie on a single line.

What we currently have...

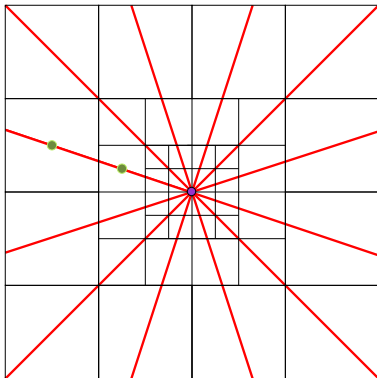
Lemma

- \mathbf{P} a set of points on a line ℓ
- ℓ - line in \mathbb{R}^d
- Coreset of size $\mathbf{O}(k/\epsilon)$ for \mathbf{P}
- For \mathbf{k} -median clustering.

Problem

General point-set do not lie on a single line.

Snapping to lines instead of grid



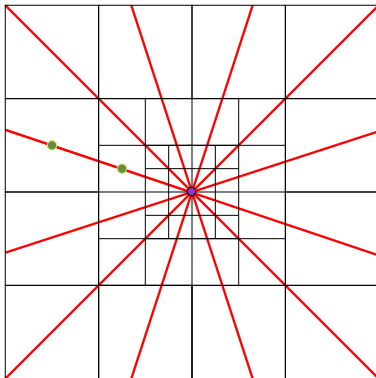
of grid cell centers: $O(\epsilon^{-d} \log n)$.

of lines: $O(1/\epsilon^d)$.

Conclusion

We can snap points to a small number of lines.

Snapping to lines instead of grid



of grid cell centers: $O(\epsilon^{-d} \log n)$.

of lines: $O(1/\epsilon^d)$.

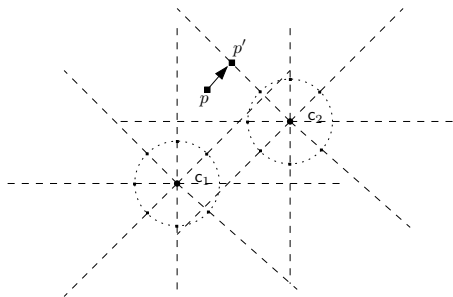
Conclusion

We can snap points to a small number of lines.

Extending to Higher Dimensions

- comp. $\mathbf{C} = \{c_1, \dots, c_k\}$ s.t. $\nu_{\mathbf{C}}(\mathbf{P}) = O(\nu_{\text{opt}}(\mathbf{P}, k))$
- construct a fan of lines around points of \mathbf{C}
- snap points of \mathbf{P} to the closest line
- build coreset for each set on a line.

Correctness



- Number of lines is $\mathbf{O}\left(\frac{k}{\epsilon^{d-1}}\right) \Rightarrow$ coreset size is $\mathbf{O}\left(\frac{k^2}{\epsilon^d}\right)$
- snapping error is bounded by $\frac{\epsilon}{3}\nu_{\text{opt}}(\mathbf{P}, \mathbf{k})$
- error by construction $\frac{\epsilon}{3}\nu_{\text{opt}}(\mathbf{P}', \mathbf{k})$ (\mathbf{P}' are snapped points)
- Total error $\leq \epsilon\nu_{\text{opt}}(\mathbf{P}, \mathbf{k})$

k-means: coresets construction

- Similar scheme works for *k*-means.
- Analysis slightly simpler.
- \Rightarrow Coreset size $O(k^3/\epsilon^{d+1})$
- Full details in the paper.

k -means: results

Definition (Centroid Set)

$\mathcal{D} \subseteq \mathbb{R}^d$: (k, ε) -approximate centroid set for \mathbf{P} , if $\exists \mathbf{C} \subseteq \mathcal{D}$ such that $\mu_{\mathbf{C}}(\mathbf{P}) \leq (1 + \varepsilon)\mu_{\text{opt}}(\mathbf{P}, k)$

- Matoušek [**Matoušek, 2000**]: \exists an ε -approx. centroid set.
Size $\mathbf{O}(n\varepsilon^{-d} \log(1/\varepsilon))$
- [**Effros and Schulman, 2003**]
Centroid set of size $\mathbf{O}(\varepsilon^{-d-1}(k^4 + k^2\varepsilon^{-2}))$.

New Result:





A (k, ε) -centroid set for k -means clustering with size $\mathbf{O}((k^3/\varepsilon^{2d+1}) \log(1/\varepsilon))$.





Fast running time...

Minor improvement in RT for k -means clustering (over [**Har-Peled and Mazumdar, 2004**])

Open Problems

- Improve running time of approx. k -means clustering?
- FPTAS for k -median and k -means (in both k and $1/\epsilon$)?
- A coresset with only polynomial dependency on the dimension and no dependency on n ?
Relevant results [Bădoiu et al., 2002, Chen, 2004].
- Improve dependency on k and ϵ in coresets size?

-  Arora, S., Raghavan, P., and Rao, S. (1998).
Approximation schemes for Euclidean k -median and related problems.
In *Proc. 30th Annu. ACM Sympos. Theory Comput.*, pages 106–113.
-  Bădoiu, M., Har-Peled, S., and Indyk, P. (2002).
Approximate clustering via coresets.
In *Proc. 34th Annu. ACM Sympos. Theory Comput.*, pages 250–257.
-  Chen, K. (2004).
Clustering algorithms using adaptive sampling.
manuscript.
-  de la Vega, W. F., Karpinski, M., Kenyon, C., and Rabani, Y. (2003).
Approximation schemes for clustering problems.
In *Proc. 35th Annu. ACM Sympos. Theory Comput.*, pages 50–58.

-  Effros, M. and Schulman, L. J. (2003).
Deterministic clustering with data nets.
Technical Report TR04-085, Elec. Colloq. Comp. Complexity.
<http://www.eccc.uni-trier.de/eccc-reports/2004/TR04-085/>.
-  Har-Peled, S. and Mazumdar, S. (2004).
Coresets for k -means and k -median clustering and their applications.
In *Proc. 36th Annu. ACM Sympos. Theory Comput.*, pages 291–300.
-  Kolliopoulos, S. G. and Rao, S. (1999).
A nearly linear-time approximation scheme for the euclidean k -median problem.
In *Proc. 7th Annu. European Sympos. Algorithms*, pages 378–389.
-  Kumar, A., Sabharwal, Y., and Sen, S. (2004).

Linear time algorithms for clustering problems in any dimension.

[manuscript](#).



[Matoušek, J. \(2000\).](#)

On approximate geometric **k**-clustering.

Discrete Comput. Geom., 24:61–84.