

Tackling the Challenges of Big Data

Big Data Analytics

Daniela Rus

Professor

Director of CSAIL

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data

Big Data Analytics

Data Compression

Daniela Rus

Professor

Director of CSAIL

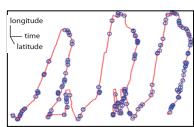
Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology



Learning Big Data Patterns from Tiny Core-sets



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Data Challenge

- 2.5 quintillion ($2.5 * 10^{18}$) bytes per day in 2012 (source: ibm.com)
- Per capita information capacity ~ doubled every 40 months since 1980
- Cameras, phones, sensors enable life-logging capabilities

Photos removed due to copyright restrictions.
We are sorry for any inconvenience this may cause.



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Data
Data
Data
Data
= Data
Data
Data

Photos removed due to
copyright restrictions.
We are sorry for any
inconvenience this may cause.



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



How much data?



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology





1 GPS Packet
= 100 bytes
(latitude, longitude, time)

MIT PROFESSIONAL EDUCATION Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology MIT CLASS

1 GPS Packet
= 100 bytes
every 10 seconds

MIT PROFESSIONAL EDUCATION Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology MIT CLASS



~ 0.4 Mb / hour
or
~100Mb / day

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

MIT PROFESSIONAL EDUCATION

DATA SCIENCE



~0.1 Gb / day
per device

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

MIT PROFESSIONAL EDUCATION

DATA SCIENCE



~300 million
smart phones
sold in 2010

<http://mobithinking.com/mobile-marketing-tools/latest-mobile-stats>

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

MIT PROFESSIONAL EDUCATION

DATA SCIENCE



For
100 million devices

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

MIT PROFESSIONAL EDUCATION

DATA
SCIENCE



For
100 million devices

~ 10 petabytes
per day

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

MIT PROFESSIONAL EDUCATION

DATA
SCIENCE



~ 10 thousand
terabytes
per day

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

MIT PROFESSIONAL EDUCATION

DATA
SCIENCE

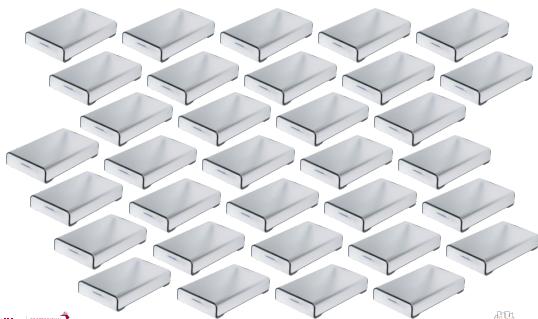
2 terabytes each



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



x5000 / day



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



That's a lot of data.



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Example: Activities

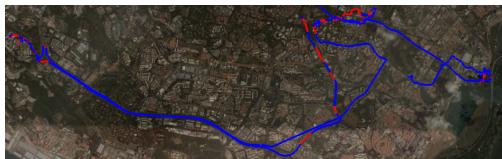
Photo removed due to copyright restrictions.
We are sorry for any inconvenience this may cause.



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Travel Mode



Blue: Predicted as if on Wheels

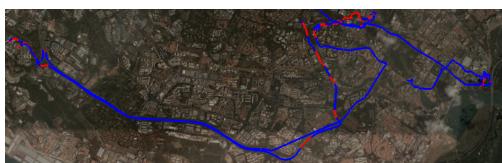
Red: Predicted as if on Foot



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Travel Mode

Recall: fraction of relevant instances classified correctly
Precision: fraction of instances correctly classified



On Foot Prediction:
93% Recall, 86% Precision



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Travel Mode

Recall: fraction of relevant instances classified correctly
Precision: fraction of instances correctly classified

**On Wheels Prediction:
96% Recall, 98% Precision**

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Example: Activities

The diagram illustrates a complex network of activities and locations. It shows a central junction with various paths leading to different locations. Each location has a unique icon and is connected to the junction by a path labeled with letters (a-h). A person is shown walking from Home to Work. The diagram is enclosed in a black-bordered box.

The diagram illustrates the process of semantic compression. It starts with a large box labeled "Time | latitude longitude" containing a grid of data points. An arrow points from this box to a smaller box labeled "Time | mode | location". Another arrow points from this second box to a third box labeled "Time | activity".

Time	Mode	Location
8:44:57	1.295783	103.7816
8:44:59	1.295785	103.7816
8:45:00	1.295782	103.7816
8:45:01	1.295782	103.7816
8:45:04	1.29579	103.7817
8:45:05	1.295782	103.7817
8:45:08	1.295915	103.7818
8:45:09	1.29598	103.7819
8:45:10	1.296015	103.7819
8:45:11	1.296057	103.782
...

Time	Activity
8:44:57	Wheel Central Ferry, Singapore
8:44:59	Wheel This Street Rd, Singapore
8:45:00	Wheel 1 S Beach Rd, Singapore 199507
8:45:01	Wheel 90A Harvey Rd, Singapore 169642
8:45:04	Foot 2 Petaling St, Singapore 050002
8:45:05	Foot Petaling St, Singapore 050002
8:45:08	Foot 212 Orchard Rd, Singapore 238832
8:45:09	Wheel 15 George Rd, Singapore 259823
8:45:10	Wheel 15 George Rd, Singapore 259823
8:45:11	Wheel Real World Council, Singapore 259823
...	...

Challenges

- Storing data on smart phone is expensive
- Transmission data is expensive
- Hard to interpret raw data
- Dynamic real-time streaming data



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Outline

- Motivation
- Coreset definition, computation model
- Coresets for k-means
- Use case: life logging systems
- Coreset for k-segments
- From coresets to text



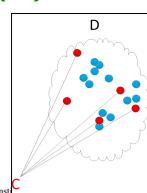
Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Challenge:

Find RIGHT data from Big Data

Given data D and Algorithm A with $A(D)$ intractable, can we efficiently reduce D to C so that $A(C)$ fast and $A(C) \sim A(D)$?

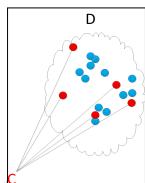


Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Challenge:**Find RIGHT data from Big Data**

Given data D and Algorithm A with $A(D)$ intractable, can we efficiently reduce D to C so that $A(C)$ fast and $A(C) \sim A(D)$?

C should be fast to compute
 $A(C) \sim A(D)$ should be provable



MIT PROFESSIONAL EDUCATION

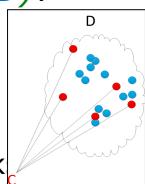
Tackling the Challenges of Big Data © 2014 Massachusetts Inst

Challenge:**Find RIGHT data from Big Data**

Given data D and Algorithm A with $A(D)$ intractable, can we efficiently reduce D to C so that $A(C)$ fast and $A(C) \sim A(D)$?

Provable guarantees for:

- Generalization
- Debugging
- Removing bottleneck



MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data © 2014 Massachusetts Inst

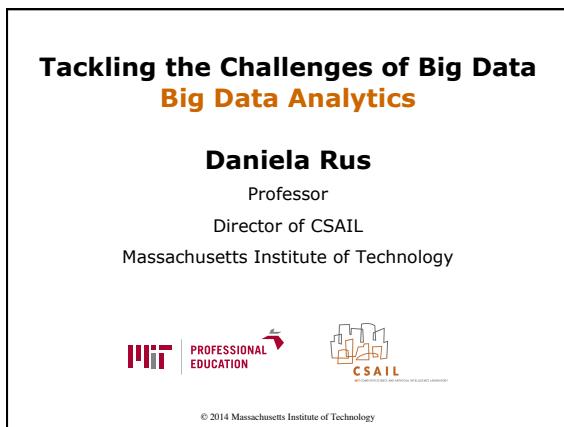
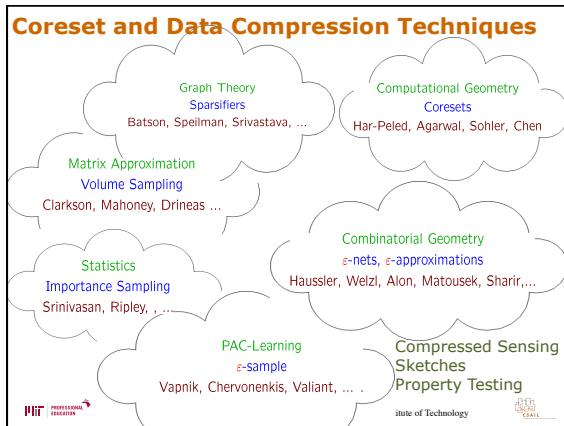
Big Data Computation Model

- Input: Infinite stream of vectors
- n vectors seen so far
- $\sim \log n$ memory
- M processors
- $\sim \log(n) / M$ insertion time per point

MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology





Tackling the Challenges of Big Data

Big Data Analytics

Data Compression

Daniela Rus

Professor
Director of CSAIL
Massachusetts Institute of Technology

 © 2014 Massachusetts Institute of Technology 

References for Compression

- k-Means [Feldman, Langberg, STOC'11]
- k-Segments [Feldman, Sung, Rus, ACM-GIS'12]
- Text Mining (LSA/PCA) [Feldman, Sohler, SODA'13]
- k-Lines [Feldman, Fiat, Sharir, FOCS'09]
- Mixture of Gaussians [Feldman, Krause, NIPS'11]
- Google Pagerank [Feldman, Yahoo! Research]
- Image Compression [Feldman, Sochen, J. of Math. Image & Vision]
- Video Compression [Feldman, Newman, submitted]

 Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology 

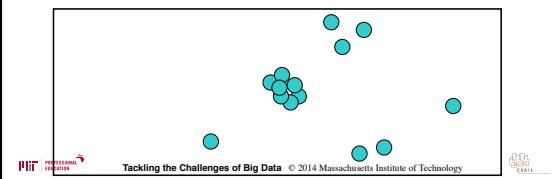
Outline

- Motivation
- Coreset definition, computation model
- Coresets for k-means
- Use case: life logging systems
- Coreset for k-segments
- From coresets to text

 Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology 

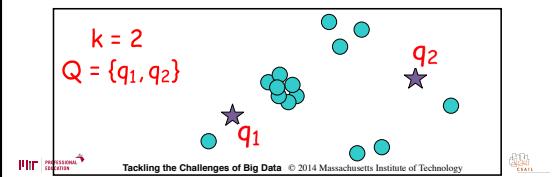
Coreset Example: K-Median Queries

- Input: $P \subseteq \mathbb{R}^d$



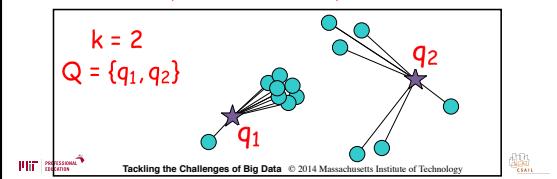
K-Median Queries

- Input: $P \subseteq R^d$
 - Query: A set Q of k points



k-Median Queries

- Input: $P \subseteq \mathbb{R}^d$
 - Query: A set Q of k points
 - Output: $\sum_{p \in P} \text{dist}(p, Q) = \sum_{p \in P} \min_{q \in Q} \|p - q\|$



(k, ϵ)-Median Coreset

Answer k -median queries in **sub-linear time**

Key Idea: Replace many points by one weighted representative:

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

(k, ϵ)-Median Coreset

Answer k -median queries in **sub-linear time**

Key Idea: Replace many points by one weighted representative:

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Definition: (k, ϵ)-Median Coreset

C is a (k, ϵ) -coreset for P , if $\forall Q, |Q| = k$:

$$\sum_{p \in P} \text{dist}(p, Q) \sim \sum_{c \in C} w(c) \cdot \text{dist}(c, Q)$$

Multiplicative error $\leq 1 + \epsilon$

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

(k, ϵ)-Median Coreset

$$\sum_{p \in P} \text{dist}(p, Q) \sim \sum_{c \in C} w(c) \cdot \text{dist}(c, Q)$$

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Can we use Naïve Uniform Sampling?

$\bullet = x_i \in \mathbb{R}^d$

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Naïve Uniform Sampling

$\bullet = x_i \in \mathbb{R}^d$

Sample a set U of m points uniformly

Small cluster is missed

← High variance

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Sampling Distribution

Bias sampling towards small clusters

Sampling distribution

Weights

How to define importance?

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

How to Define Distribution? Weights

Far: high prob, low weight
Near: low prob, high weight

Sampling distribution

Weights

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Creating a Sampling Distribution

Points in sparse cells get more mass
Points far from centers get more mass

$q(x)$

Find a good first guess
Sample wrt distance + $1/\| \text{points in cluster} \|$

Sampling distribution q

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Theorem: general claim for coresets

[Feldman, Langberg, STOC'11]

Let P, Q and $\text{dist}: P \times Q \rightarrow R^+$.

A sample $C \subseteq P$ from the distribution

$$\text{sensitivity}(p) = \max_{q \in Q} \frac{\text{dist}(p, q)}{\sum_{p'} \text{dist}(p', q)}$$

is a coreset if $|C| \geq \frac{\text{dimension of } Q}{\epsilon^2} \cdot \sum_p \text{sensitivity}(p)$



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Coreset for k-means

[Feldman, Sohler, Monemizadeh, SoCG'07]

The coreset for k-means can be computed by choosing points from the distribution:

$$\text{sensitivity}(p) = \frac{\text{dist}(p, q^*)}{\sum_{p'} \text{dist}(p', q^*)} + \frac{1}{n_p}$$

q^* = k-means of P

n_p = number of points in the cluster of p

$$|C| = \frac{k \cdot d}{\epsilon^2}$$

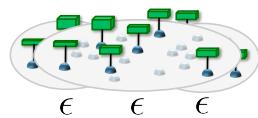


Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Composition of Coresets

Merge The union of two (k, ϵ) -coresets is a (k, ϵ) -coreset.

- [c.f. Har-Peled, Mazumdar 04]

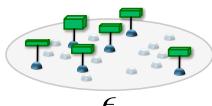


Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Composition of Coresets

Merge The union of two (k, ϵ) -coresets is a (k, ϵ) -coreset.

Compress A (k, δ) -coreset of a (k, ϵ) -coreset is a $(k, \epsilon + \delta + \epsilon\delta)$ -coreset

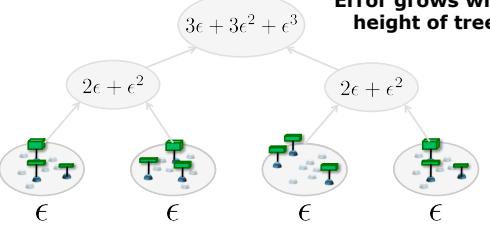


$$\epsilon + \delta + \epsilon\delta$$

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coresets on Streams

Error grows with height of tree



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coresets in Parallel

Diagram illustrating coresets in parallel, showing a tree structure where each node is a coreset. The root node has an error of $3\epsilon + 3\epsilon^2 + \epsilon^3$. Its children have errors of $2\epsilon + \epsilon^2$. The leaf nodes have errors of ϵ .

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data Big Data Analytics

Data Compression

THANK YOU



© 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data Big Data Analytics

Daniela Rus

Professor

Director of CSAIL

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data Big Data Analytics

Data Compression

Daniela Rus

Professor

Director of CSAIL

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology



Outline

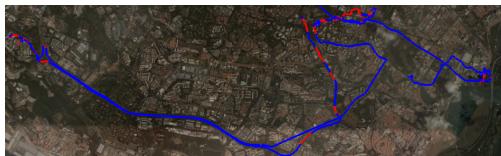
- Motivation
- Coreset definition, computation model
- Coresets for k-means
- Use case: life logging systems
- Coreset for k-segments
- From coresets to text



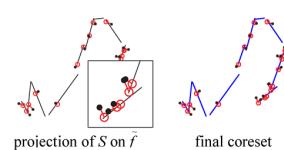
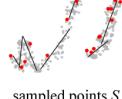
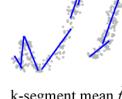
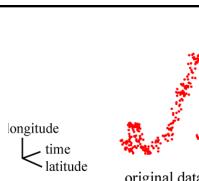
Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Coresets for Travel: k-segments



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



final coreset

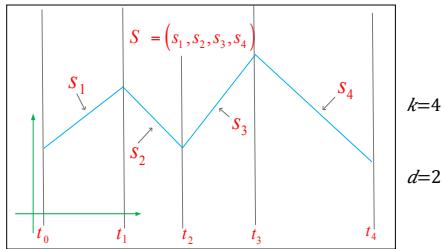


Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Definition: k-spline

A k-spline is a sequence of k connected segments in \mathbb{R}^d

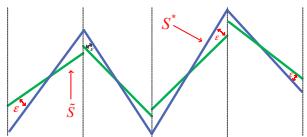


Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



ϵ -approximation for $\text{Opt}(P, k)$

A set \tilde{S} of k segments is an ϵ -approximation for $\text{Opt}(P, k)$ if rotating every segment of \tilde{S} by angle of at most $\epsilon > 0$ yields an optimal k -spline S^* of P

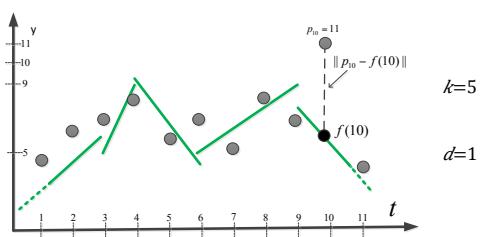


Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



k-segment Mean

The k-segment that minimizes the fitting costs from points to a d-dimensional signal



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



k-segment Queries

Input: d -dimensional signal P over time

MIT PROFESSIONAL EDUCATION Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

k-segment Queries

Input: d -dimensional signal P over time
Query: k segments over time

MIT PROFESSIONAL EDUCATION Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

k-segment Queries

Input: d -dimensional signal P over time
Query: k segments over time
Output: Sum of squared distances from P

MIT PROFESSIONAL EDUCATION Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset for k-segments

- A weighted set of segments C such as for every k -segment f

$$\text{cost}(P, f) \sim \text{cost}_w(C, f)$$

$\sum_t \|f(t) - p_t\|^2$ $(1 \pm \epsilon)$ $\sum_{p_t \in C} w(p_t) \cdot \|f(t) - p_t\|^2$

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset Compression Theorem

[ACM GIS'12, Feldman, Sung, and Rus]

For every discrete signal of n points in R^d there is a

- coreset of space $O(k/\epsilon^2)$
- that can be computed in the big data model

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset Construction

Input: signal of n points,
constants k, ϵ

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset Construction

Input: signal of n points,
constants k, ε

Compute k -segment mean

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset Construction

Input: signal of n points,
constants k, ε

Compute k -segment mean

Project points onto segments

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset Construction

Input: signal of n points,
constants k, ε

Compute k -segment mean

Project points onto segments

Assign probability $\sigma_p = \frac{d_p}{\sum_p d_p}$

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset Construction

Input: signal of n points,
constants k, ϵ

Compute k -segment mean

Project points onto segments

Sample $|S| \sim \frac{k}{\epsilon^2}$ points
with probability $\sigma_p = \frac{d_p}{\sum_p d_p}$

coresets

Coreset Construction

Input: signal of n points,
constants k, ϵ

Compute k -segment mean

Project points onto segments



Sample $|S| \sim \frac{k}{\epsilon^2}$ points
with probability $\sigma_p = \frac{d_p}{\sum_p d_p}$

Assign weights $\pm \frac{1}{|S|\sigma_p}$




Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Coreset Construction

Input: signal of n points,
constants k, ϵ, δ

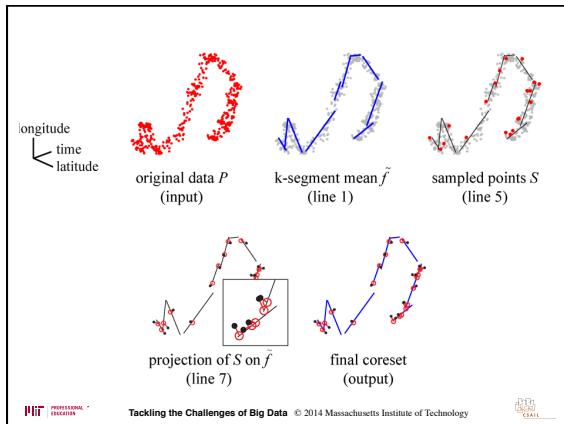
Compute k -segment mean

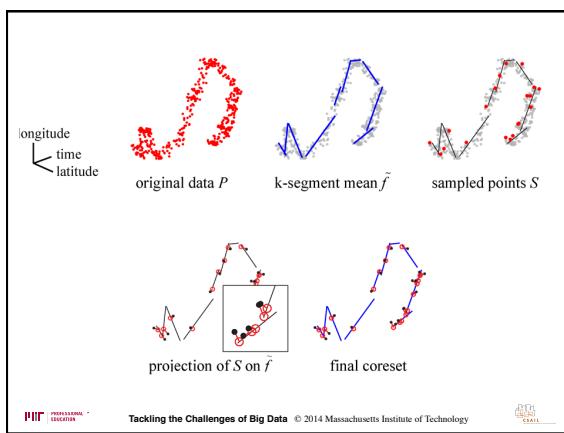
Project points onto segments

Sample $\lceil S/\epsilon \rceil$ points
with probability $\alpha p = d$

Assign weights $\pm 1/\lceil S/\alpha p \rceil$

Output: k segments and
 $|S|$ weighted points pairs







Tackling the Challenges of Big Data Big Data Analytics

Daniela Rus

Professor

Director of CSAIL

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data Big Data Analytics Data Compression

Daniela Rus

Professor

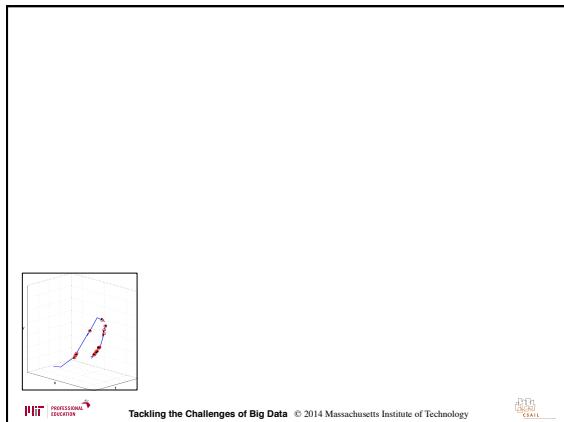
Director of CSAIL

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

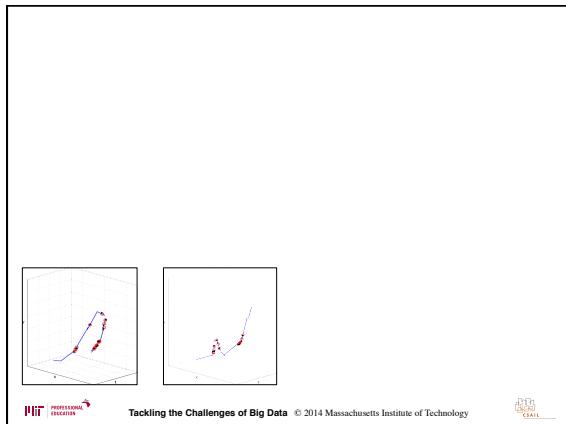


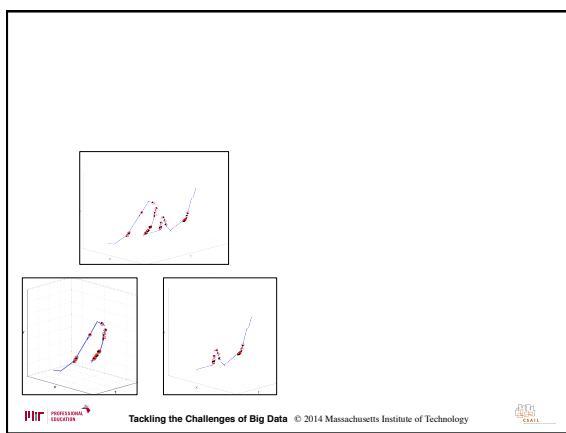


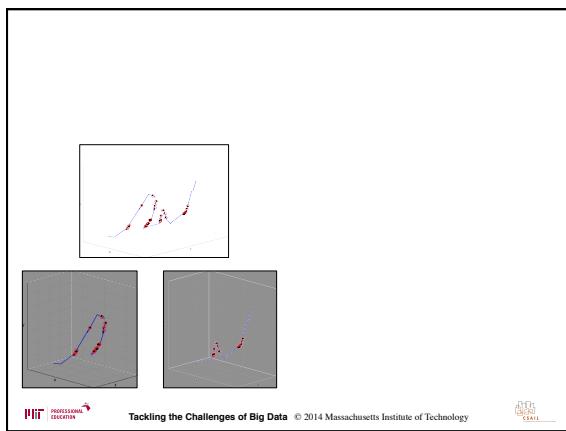
©

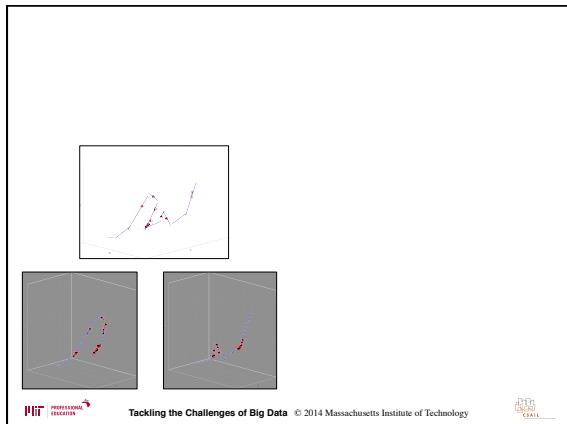
Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



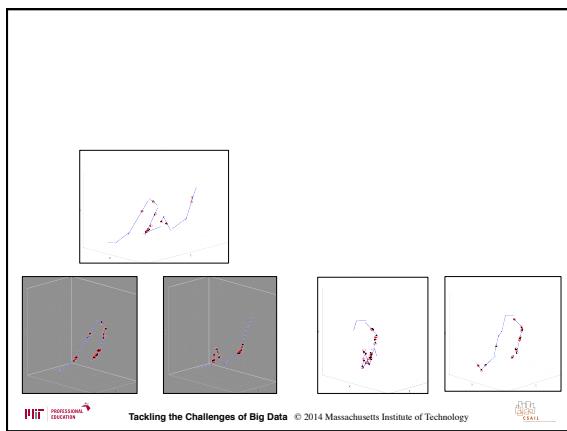


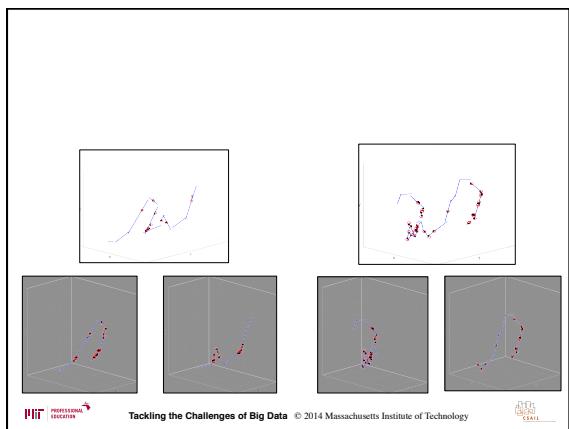


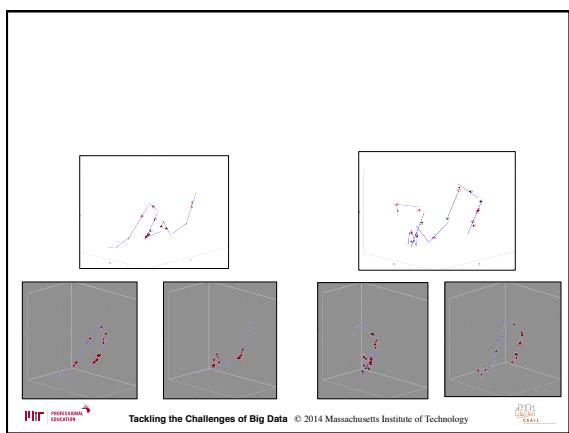


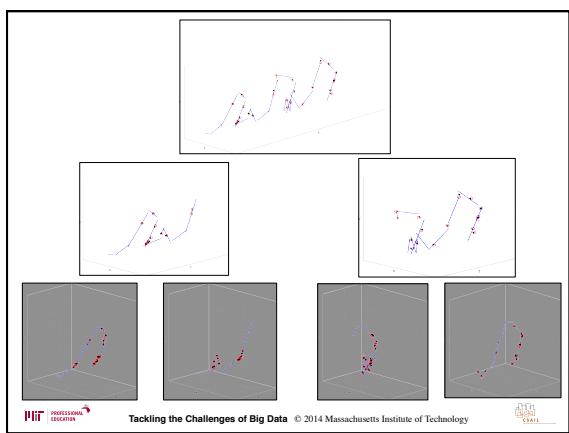


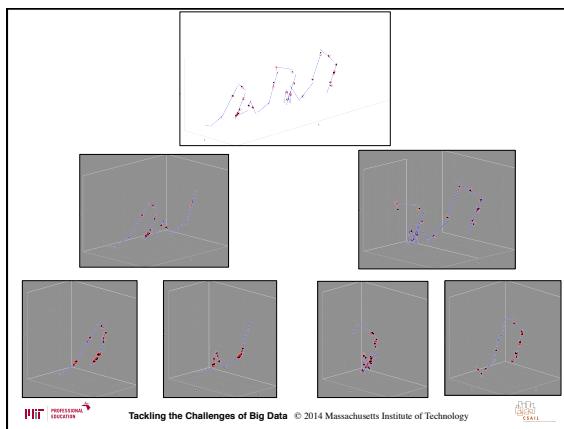


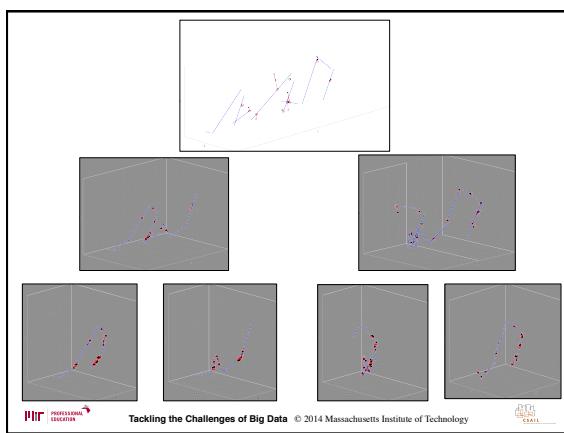


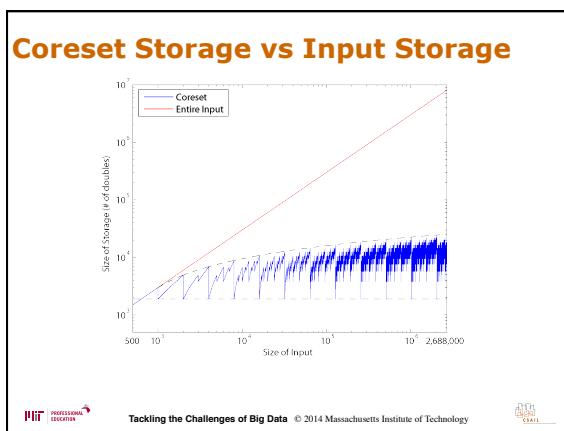












Computation

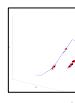
- We need to apply the coresets only on some points
- The construction time is $O(k^3)$ instead of $O(n^3)$
- If k^3 is still too long, or we don't know how to compute optimal solution
 - Use bi-criteria approximations
 - or heuristics



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



Parallel Computation



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

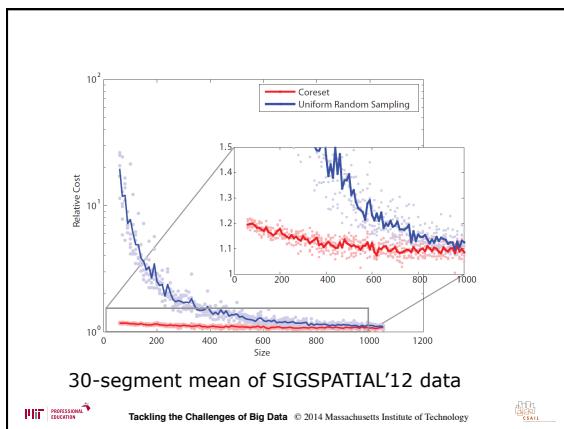
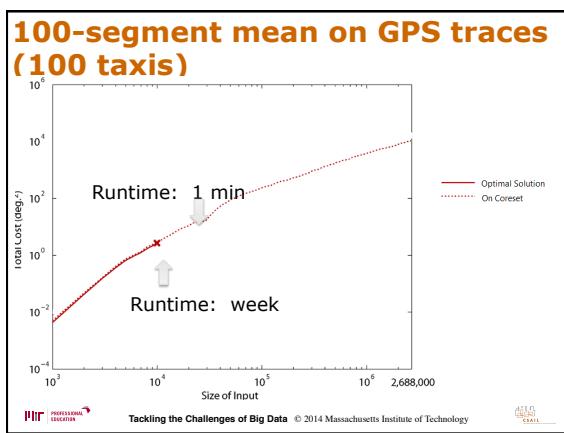
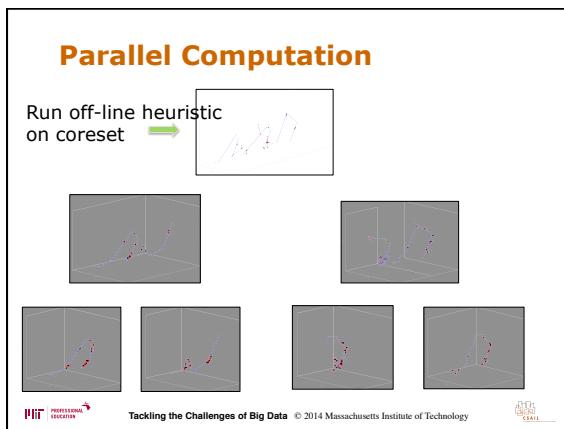


Parallel Computation



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology





Tackling the Challenges of Big Data Big Data Analytics

Data Compression

THANK YOU



© 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data Big Data Analytics

Daniela Rus

Professor

Director of CSAIL

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data Big Data Analytics

Data Compression

Daniela Rus

Professor

Director of CSAIL

Massachusetts Institute of Technology

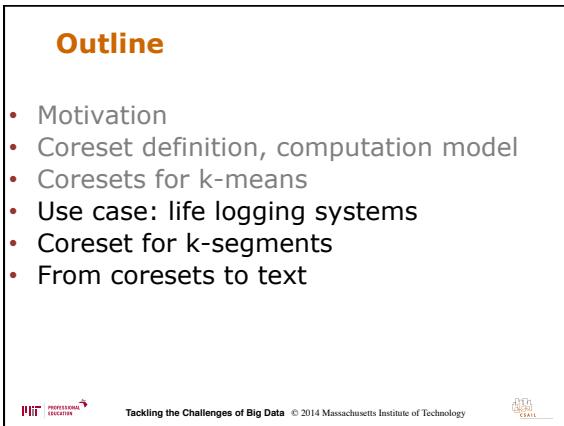


© 2014 Massachusetts Institute of Technology



Outline

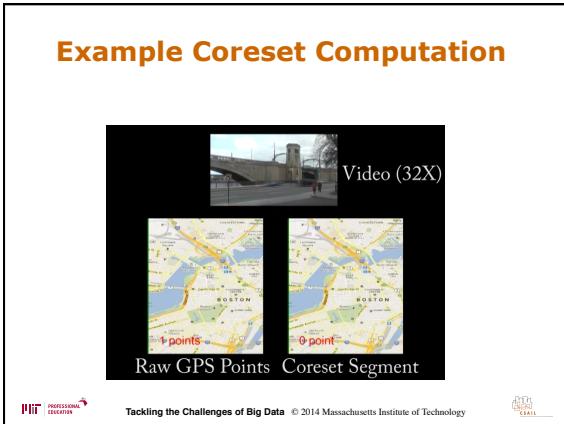
- Motivation
 - Coreset definition, computation model
 - Coresets for k-means
 - Use case: life logging systems
 - Coreset for k-segments
 - From coresets to text

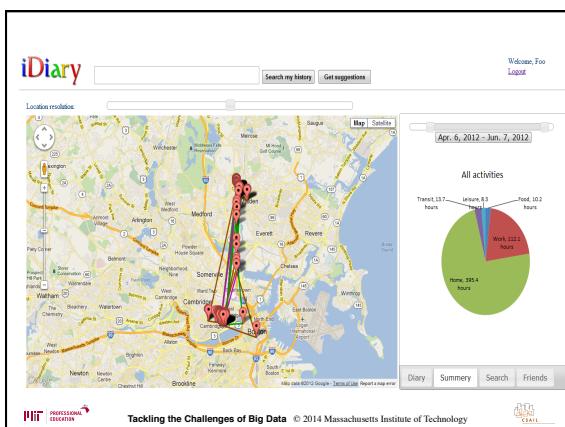
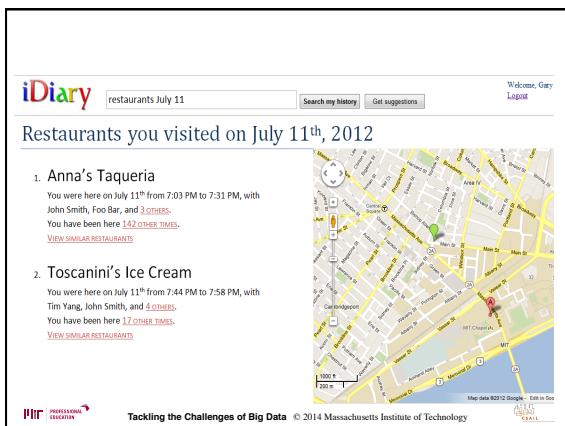


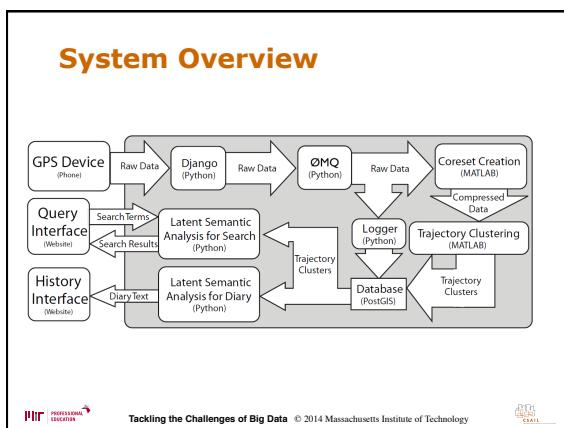
Example: Coresets for Life Logging



Example Coreset Computation







Coreset Computation

- Approximate the input GPS-points by connected segments using a k-spline
- Output the text description of the endpoints (e.g., using Google Maps)

time	latitude	longitude
8:44:57	1.295783	103.7816
8:44:59	1.295785	103.7816
8:45:00	1.295782	103.7816
8:45:02	1.295782	103.7816
8:45:04	1.295787	103.7817
8:45:05	1.295802	103.7817
8:45:06	1.295914	103.7818
8:45:09	1.29598	103.7819
8:45:10	1.296013	103.7819
8:45:11	1.296057	103.782
...

$k = 20$

$k+1 = 21$

The figure shows a 3D plot of GPS points (blue dots) and a fitted k-spline (blue line). A red arrow points from the data points to the spline. Another red arrow points from the spline to the text output on the right.

Final Endpoints:

- 1) 103.7816, 1.295783
- 2) 103.7816, 1.295785
- 3) 103.7816, 1.295782
- 4) 103.7816, 1.295782
- 5) 103.7817, 1.295787
- 6) 103.7817, 1.295802
- 7) 103.7818, 1.295914
- 8) 103.7819, 1.29598
- 9) 103.7819, 1.296013
- 10) 103.782, 1.296057

Intermediate Endpoints:

- 1) 103.7816, 1.295783
- 2) 103.7816, 1.295785
- 3) 103.7816, 1.295782
- 4) 103.7816, 1.295782
- 5) 103.7817, 1.295787
- 6) 103.7817, 1.295802
- 7) 103.7818, 1.295914
- 8) 103.7819, 1.29598
- 9) 103.7819, 1.296013
- 10) 103.782, 1.296057

Geographic Labels:

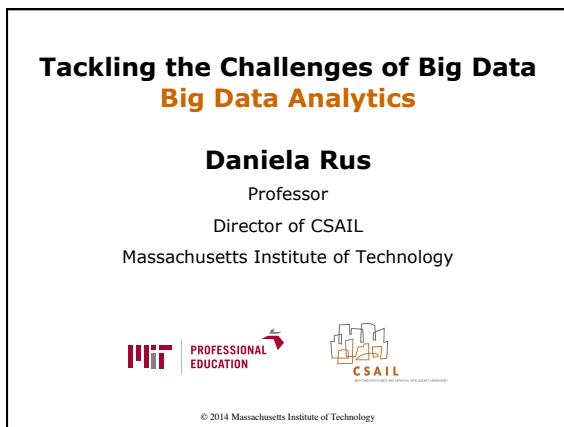
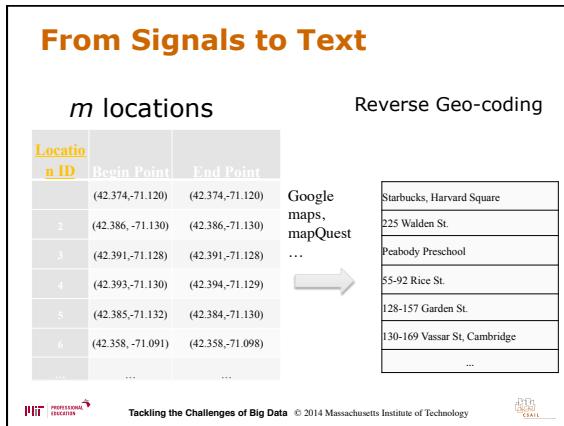
- 1) Jalan Sungai Selatan
- 2) Jalan Sungai Selatan
- 3) Jalan Sungai Selatan
- 4) Jalan Sungai Selatan
- 5) Jalan Sungai Selatan
- 6) Jalan Sungai Selatan
- 7) Jalan Sungai Selatan
- 8) Jalan Sungai Selatan
- 9) Jalan Sungai Selatan
- 10) Jalan Sungai Selatan

Map Labels:

- 1) Jalan Sungai Selatan
- 2) Jalan Sungai Selatan
- 3) Jalan Sungai Selatan
- 4) Jalan Sungai Selatan
- 5) Jalan Sungai Selatan
- 6) Jalan Sungai Selatan
- 7) Jalan Sungai Selatan
- 8) Jalan Sungai Selatan
- 9) Jalan Sungai Selatan
- 10) Jalan Sungai Selatan

Geotag Labels:

- 1) Jalan Sungai Selatan
- 2) Jalan Sungai Selatan
- 3) Jalan Sungai Selatan
- 4) Jalan Sungai Selatan
- 5) Jalan Sungai Selatan
- 6) Jalan Sungai Selatan
- 7) Jalan Sungai Selatan
- 8) Jalan Sungai Selatan
- 9) Jalan Sungai Selatan
- 10) Jalan Sungai Selatan



Tackling the Challenges of Big Data

Big Data Analytics

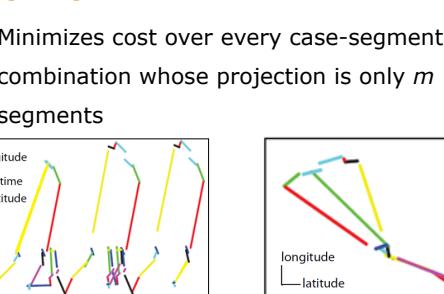
Data Compression

Outline

- Motivation
- Coreset definition, computation model
- Coresets for k-means
- Use case: life logging systems
- Coreset for k-segments
- From coresets to text

(k,m) Formulation

Minimizes cost over every case-segment combination whose projection is only m segments



Travel Pattern Input

- GPS-point = (latitude, longitude, time)

latitude	longitude	time
1.295783	103.7816	8:44:57
1.295783	103.7816	8:44:59
1.295783	103.7816	8:45:00
1.295783	103.7816	8:45:01
1.295776	103.7817	8:45:04
1.295802	103.7817	8:45:05
1.295802	103.7818	8:45:08
1.295806	103.7819	8:45:09
1.296013	103.7819	8:45:10
1.296057	103.782	8:45:11

longitude

time

latitude

longitude

latitude

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

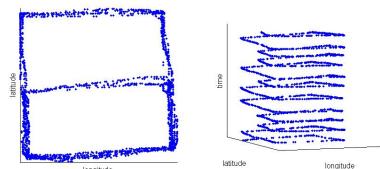
The figure illustrates travel patterns. On the left, two sets of trajectories are shown as colored lines connecting points on a grid. The top set is labeled 'k trajectories' and the bottom set is labeled 'm locations'. On the right, a single set of trajectories is shown as colored lines connecting points on a grid.

Begin time	End time	Location ID	Speed
8:45-57	8:45-57	c	30
8:45-59	8:51-59	d	24
8:55-60	8:54-60	g	24
8:55-61	8:55-61	q	11
8:55-57	8:57-57	r	120
8:55-56	8:59-57	m	55
...	65

Location ID	Begin Point	End Point
a	(42.374, -71.120)	(42.374, -71.120)
b	(42.386, -71.130)	(42.386, -71.130)
c	(42.391, -71.128)	(42.391, -71.128)
d	(42.393, -71.130)	(42.394, -71.129)

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

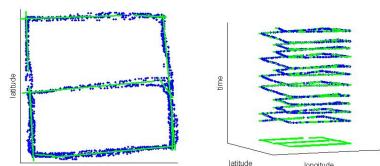
A GPS signal: a sequence of geographic points progressing through time
 Derived from an entity's repeated traversal of a set of distinct geographic paths



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



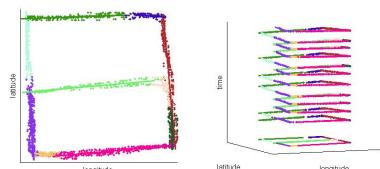
But where are those paths?
 The human eye can make a good guess
 Can an algorithm do the same?



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



A (k, m) -segment trajectory:
 • Partition the signal into k sections
 • Group the sections into m clusters
 • Assign a line segment to each cluster

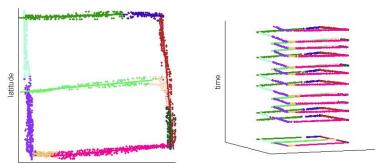


Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



The (k, m) -segment mean:

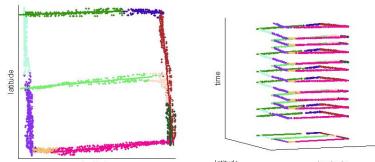
- A (k, m) -segment trajectory
 - Minimizes the point signal's fitting cost



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Our (k, m) -segment mean algorithm:

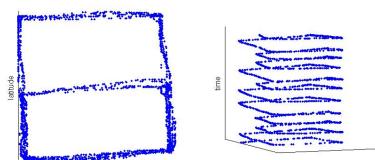
- Approximates the true mean
 - Uses EM to find local minimum of fitting cost
 - Efficient in runtime and space



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

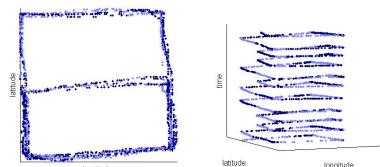
Start with the GPS point signal

No external environment data (e.g. road maps)



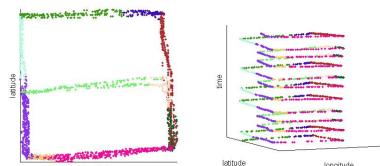
Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Partition the signal into k distinct sections
Use Ramer-Douglas-Peucker for a good first guess (it targets high-cost points)



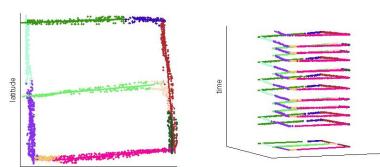
Tackling the Challenges of Big Data. © 2014 Massachusetts Institute of Technology.

Group k sections into m clusters (colors)
Use k-means for a good first guess (clustering according to endpoint locations)



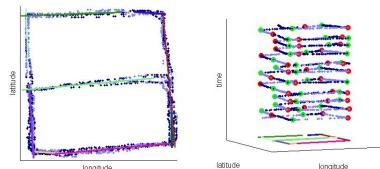
Tackling the Challenges of Big Data. © 2014 Massachusetts Institute of Technology

Calculate the segment mean of each cluster
Using linear algebra makes it exact (non-approximate)



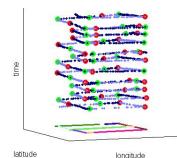
longitude longitude longitude

Hold fixed the odd partition boundaries (red)
 Given the m segments, optimally adjust the even
 boundaries (green) and the section cluster
 assignments



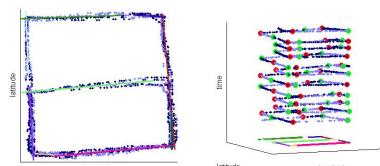
MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



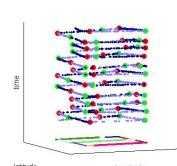
MIT PROFESSIONAL EDUCATION

Hold fixed the even partition boundaries (red)
 Given the m segments, optimally adjust the odd
 boundaries (green) and the section cluster assignments



MIT PROFESSIONAL EDUCATION

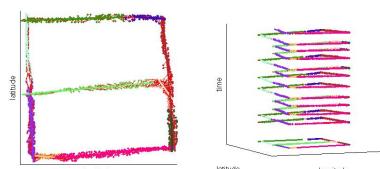
Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology



MIT PROFESSIONAL EDUCATION

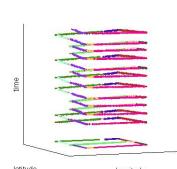
With the updated partition and clusters, recalculate each
 cluster's segment mean, and the total fitting cost (red
 lines)

If cost is reduced repeat; else terminate



MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

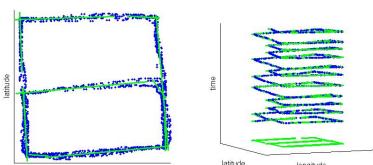


MIT PROFESSIONAL EDUCATION

Upon termination, the locally optimal (k, m) -segment trajectory has been found

Line segments give the map

Clustering gives trajectory on map



Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Two domains of algorithmic effectiveness

- Accuracy (minimal error)
 - Speed (minimal runtime)

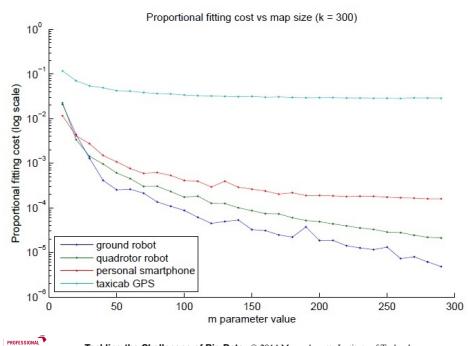
Measure of accuracy: proportional fitting cost

- Average squared error distance per input point, relative to points' variance (squared error from mean)
 - Independent of input's set size and geographic magnitude
 - In range $[0, 1]$ for exact (k, m) -segment mean

Measure of speed: runtime per point

- Average runtime per point
 - Independent of input's set size and geographic magnitude

Accuracy across four sample data sets

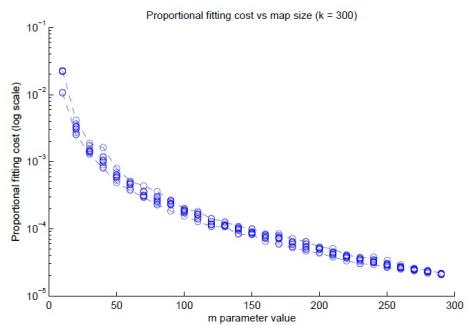


MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

DATA
SCIENCE
CLOUD

Accuracy across ten runs of a data set

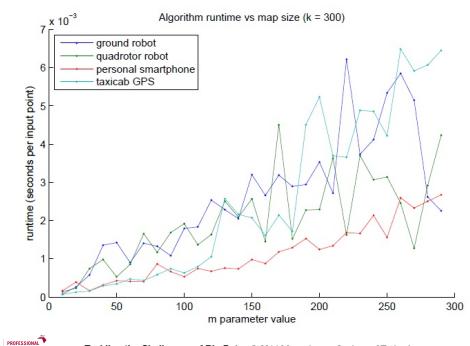


MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

DATA
SCIENCE
CLOUD

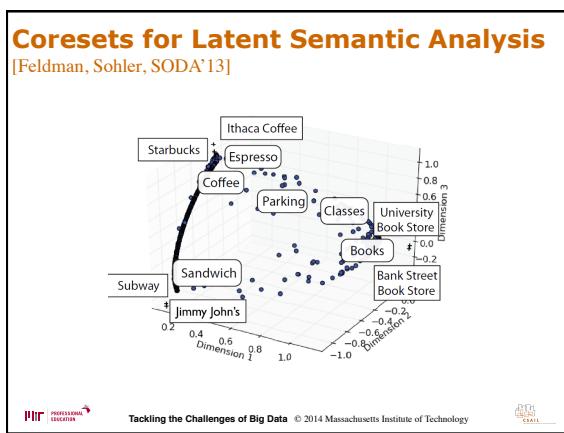
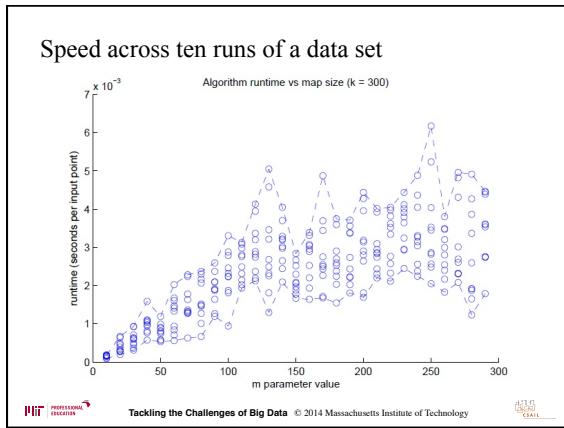
Speed across four sample data sets



MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

DATA
SCIENCE
CLOUD



- ## Wrapping Up
- Instead of Speeding Algorithms, Find Right Data
 - Coresets provide semantic compression for data
 - Sample coreset algorithms
 - Applications for personal travel management, transportation
- Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

Tackling the Challenges of Big Data

Big Data Analytics

Data Compression

THANK YOU



© 2014 Massachusetts Institute of Technology
