Prioritized Geographic Search

SPP subproject "Lightweight Acquisition and Large-Scale Mining of Trajectory Data"



Prof. Stefan Funke

Institut f. Formale Methoden der Informatik Abteilung Algorithmik



University of Stuttgart

Germany



Given: Set S of n points in \mathbb{R}^2



Given: Set S of n points in \mathbb{R}^2

Goal: Construct data structure \mathcal{D} such that *range queries* specified by $[x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$ can be answered efficiently.



Given: Set S of n points in \mathbb{R}^2

Goal: Construct data structure \mathcal{D} such that *range queries* specified by $[x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$ can be answered efficiently.

Examples:

Output within a given area:

- all Greek restaurants
- locations of all mailboxes



Given: Set S of n points in \mathbb{R}^2

Goal: Construct data structure \mathcal{D} such that *range queries* specified by $[x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$ can be answered efficiently.

Examples:

Output within a given area:

- all Greek restaurants
- locations of all mailboxes

We care about:

- \bullet time Q(n) to answer a single query
- \bullet space S(n) of the data structure ${\cal D}$
- $\bullet \mbox{ time } P(n) \mbox{ to construct } \mathcal{D}$

Basic Geographic Range Search – Challenges

Consider the query "all towns within the query rectangle"



Basic Geographic Range Search – Challenges

Consider the query "all towns within the query rectangle"



The result set is likely to contain several thousands of items ...

Basic Geographic Range Search – Challenges

Consider the query "all towns within the query rectangle"



The result set is likely to contain several thousands of items ... Depending on the purpose of the query, it is pretty hard to make sense of this result set ...

Web Search:

- Query for "geographic information systems" yields
 ≈ 46 Million results / webpages related to the search term
- yet, we can make sense of the outcome as Google *prioritizes* the outcome in a sensible manner
- good prioritization and prioritized querying main reason for Google prevailing over Yahoo, Altavista, ...

Web Search:

- Query for "geographic information systems" yields
 ≈ 46 Million results / webpages related to the search term
- yet, we can make sense of the outcome as Google *prioritizes* the outcome in a sensible manner
- good prioritization and prioritized querying main reason for Google prevailing over Yahoo, Altavista, ...



	Feedback
Why GIS is important for you?	\sim
What is a GIS ESRI?	\sim
What is remote sensing and geographic information system?	\sim
What GIS is all about?	\sim
People also ask	

Geographic information system - Wikipedia

https://en.wikipedia.org/wiki/Geographic_information_system A geographic information system (GIS) is a system designed to capture, store, manipulate, analyze, manage, and present spatial or geographic data. Spatial data infrastructure · Geographic information science · Canada Geographic ...

What is GIS? | The Power of Mapping - Esri www.esri.com/what-is-gis *

A geographic information system (GIS) lets us visualize, question, analyze and interpret data to understand relationships, patterns & trends. Learn more about ...

GIS (geographic information system) - National Geographic Society www.nationalgeographic.org/encyclopedia/geographic-information-system-gis/ •

Jun 21, 2017 - A geographic information system (GIS) is a computer system for capturing, storing, checking, and displaying data related to positions on Earth's surface. By relating seemingly unrelated data, CIS can belie individuals and arrangizitizes better understand apartial potters and

Geographic Search:

- the OpenStreetMap planet data set comprises more than 3 billion items
- an individual query might return hundreds of thousands of items
- again: *prioritization* key for reasonable use and interpretation of results



Addresses and postcodes are approximate © OpenStreetMap contributors

@shop:mall

Australia (Wine)

Store)

(Organic)

Geographic Search:

- the OpenStreetMap planet data set comprises more than 3 billion items
- an individual query might return hundreds of thousands of items
- again: *prioritization* key for reasonable use and interpretation of results

Query: @shop:mall

- \approx 2,700 in Germany
- prioritization according to size seems very reasonable



Addresses and postcodes are approximate

© OpenStreetMap contributors

When drawing a zoomed-out view of a map, only the most important (largest?) cities should be drawn ...



When drawing a zoomed-out view of a map, only the most important (largest?) cities should be drawn

When zooming-in, also medium-sized towns should be on the rendering ...



Besigheim

Bietigheim

Bissingen

Korntal

Ludwigsburg

Kornwestheim

Stuttgart

Metzingen

Reutlingen Bad Urach

When drawing a zoomed-out view of a map, only the most important (largest?) cities should be drawn

When zooming-in, also medium-sized towns should be on the rendering ...

Vaihingen

an der Enz

ler Stadt

Leonberg

Sindelfingen

Holzgerlingen

Rottenhurg

Tübingen

Danmark Hamburg Berlin Nederland Polska Deutschland Česko Paris Slovens München บมากันธ Magyarorsz Hessen ◎ Zagreb Koblenz Frankfurt am Würzburg Main Mannheim Nürnberg rücken Karlsruhe Ingolstadt Badenbourg Augsburg Württemberg München Murrh Freiburg n Breisgau Backnang Winterthur Winnenden Waiblingen Schorndor Wendlingen am Neckar Nürtingen Owen

The further we zoom-in, the more smaller villages we want to see on the map

When drawing a zoomed-out view of a map, only the most important (largest?) cities should be drawn

> When zooming-in, also medium-sized towns should be on the rendering ...



The further we zoom-in, the more smaller villages we want to see on the map ...



(Almost) direct application of prioritized range queries!

Two-dimensional Range Queries (no priorities yet)

1-dimensional range queries:

- use tree structure
- logarithmic depth
- space consumption for n elements: S(n)=O(n)
- query time for with result size k: Q(n)=O(log n + k)



Two-dimensional Range Queries (no priorities yet)

1-dimensional range queries:

- use tree structure
- logarithmic depth
- space consumption for n elements: S(n)=O(n)
- query time for with result size k: Q(n)=O(log n + k)

Naive generalization to 2 dimensions:

- build tree on x-coordinates
- build tree on y-coordinates
- S(n)=O(n)
- query x- and y-ranges separately, return intersection
- Q(n)=?



Two-dimensional Range Queries (no priorities yet)

1-dimensional range queries:

- use tree structure
- logarithmic depth
- space consumption for n elements: S(n)=O(n)
- query time for with result size k: Q(n)=O(log n + k)

Naive generalization to 2 dimensions:

- build tree on x-coordinates
- build tree on y-coordinates
- S(n)=O(n)
- query x- and y-ranges separately, return intersection
- Q(n)=?

 $Q(n) = \Theta(n)$ for a result size of 0!



1-dimensional tree construction:

- recursively via median splitting
- construction time $P(n) = O(n \log n)$



1-dimensional tree construction:

- recursively via median splitting
- \bullet construction time $P(n) = O(n \log n)$

Generalization:

- alternate between median splitting according to x- and y-coordinate
- \bullet construction time $O(n\log n)$
- \bullet space consumption S(n)=O(n)
- how to query? query time?







1-dimensional tree construction:

- recursively via median splitting
- \bullet construction time $P(n) = O(n \log n)$

Generalization:

- alternate between median splitting according to x- and y-coordinate
- \bullet construction time $O(n\log n)$
- space consumption S(n) = O(n)
- how to query? query time?





kd-Tree

- associate with each node rectangle containing all points of the subtree
- recursively traverse rectangles with non-empty intersection with query rectangle
- query time Q(n) = ?





- associate with each node rectangle containing all points of the subtree
- recursively traverse rectangles with non-empty intersection with query rectangle
- query time Q(n) = ?





- associate with each node rectangle containing all points of the subtree
- recursively traverse rectangles with non-empty intersection with query rectangle
- query time Q(n) = ?





- associate with each node rectangle containing all points of the subtree
- recursively traverse rectangles with non-empty intersection with query rectangle
- query time Q(n) = ?





- Query time can be $Q(n) = \Omega(\sqrt{n})$ with 0 points reported :(
- but it cannot be worse $Q(n) = O(\sqrt{n} + k)$

- every point now also bears a priority
- query is of the form

 $(x_{\min}, x_{\max}, y_{\min}, y_{\max}, prio_{\min})$

- every point now also bears a priority
- query is of the form

 $(x_{\min}, x_{\max}, y_{\min}, y_{\max}, prio_{\min})$

We could treat priority as 3rd dimension (allowing also max-priorities)

- $\bullet \ {\rm size} \ S(n) = O(n)$
- construction time $P(n) = O(n \log n)$
- query time $Q(n) = O(n^{2/3} + k)$

- every point now also bears a priority
- query is of the form

 $(x_{\min}, x_{\max}, y_{\min}, y_{\max}, prio_{\min})$

We could treat priority as 3rd dimension (allowing also max-priorities)

- size S(n) = O(n)
- construction time $P(n) = O(n \log n)$
- query time $Q(n) = O(n^{2/3} + k)$

Better alternative:

- Augmented kd-prio tree:
 - before median extraction, extract max-priority point and store within node
- $\bullet \text{ space } S(n) = O(n \log n)$
- \bullet construction time $P(n) = O(n \log n)$
- \bullet query time $Q(n) = O(\sqrt{n} + k)$
- \bullet abort subtree inspection as soon as priority is below $prio_{\min}$

- every point now also bears a priority
- query is of the form

 $(x_{\min}, x_{\max}, y_{\min}, y_{\max}, prio_{\min})$

We could treat priority as 3rd dimension (allowing also max-priorities)

- size S(n) = O(n)
- construction time $P(n) = O(n \log n)$
- query time $Q(n) = O(n^{2/3} + k)$

Better alternative:

- Augmented kd-prio tree:
 - before median extraction, extract max-priority point and store within node
- $\bullet \text{ space } S(n) = O(n \log n)$
- construction time $P(n) = O(n \log n)$
- \bullet query time $Q(n) = O(\sqrt{n} + k)$
- \bullet abort subtree inspection as soon as priority is below $prio_{\min}$



- organize points in 1-dimensional search tree according *x*-coordinate
- for query $[x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$:
 - search for x_{\min} and x_{\max} in tree
 - observation: all points with matching x-coordinate 'enclosed' by search paths
- store in each internal node all points of subtree in a tree on the y-coordinates
- $\bullet \text{ space } S(n) = O(n \log n)$
- query time $Q(n) = O(log^2n + k)$



- organize points in 1-dimensional search tree according *x*-coordinate
- for query $[x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$:
 - search for x_{\min} and x_{\max} in tree
 - observation: all points with matching x-coordinate 'enclosed' by search paths
- store in each internal node all points of subtree in a tree on the y-coordinates
- $\bullet \text{ space } S(n) = O(n \log n)$
- \bullet query time $Q(n) = O(log^2n + k)$



- organize points in 1-dimensional search tree according *x*-coordinate
- for query $[x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$:
 - search for x_{\min} and x_{\max} in tree
 - observation: all points with matching x-coordinate 'enclosed' by search paths
- store in each internal node all points of subtree in a tree on the y-coordinates
- $\bullet \text{ space } S(n) = O(n \log n)$
- \bullet query time $Q(n) = O(log^2n + k)$

Priorities can be incorporated as 3rd dimension:

• $S(n) = O(n \log^2 n)$ • $Q(n) = O(log^3 n + k)$



- organize points in 1-dimensional search tree according *x*-coordinate
- for query $[x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$:
 - search for x_{\min} and x_{\max} in tree
 - observation: all points with matching x-coordinate 'enclosed' by search paths
- store in each internal node all points of subtree in a tree on the y-coordinates
- $\bullet \text{ space } S(n) = O(n \log n)$
- \bullet query time $Q(n) = O(log^2n + k)$

Priorities can be incorporated as 3rd dimension:

•
$$S(n) = O(n \log^2 n)$$

• $Q(n) = O(log^3 n + k)$





20

• A treap can answer queries of the type $[x_{\min}, x_{\max}] \times [y_{\min}, \infty]$ • treap = tree and heap at the same time (41) (22) (39) (17) (13)



• A treap can answer queries of the type $[x_{\min}, x_{\max}] \times [y_{\min}, \infty]$ • treap = tree and heap at the same time (41) (22) (39) (17) (13) (20)



Construction:

- \bullet extract max-prio point and x-median point
- recursively construct left and right subtrees
- $\bullet \text{ time } P(n) = O(n \log n)$
- $\bullet \text{ space } S(n) = O(n)$

A treap can answer queries of the type

[x_{min}, x_{max}] × [y_{min}, ∞]

treap≡ tree and heap

at the same time
⁴¹
²²
³⁹
¹⁷
¹³

Construction:

- \bullet extract max-prio point and x-median point
- recursively construct left and right subtrees
- $\bullet \text{ time } P(n) = O(n \log n)$
- $\bullet \text{ space } S(n) = O(n)$

• ignoring median points, we have a heap!



22

20

13

17

A treap can answer queries of the type

 [x_{min}, x_{max}] × [y_{min}, ∞]

 treap≡ tree and heap

 the same time

Construction:

- \bullet extract max-prio point and x-median point
- recursively construct left and right subtrees
- $\bullet \text{ time } P(n) = O(n \log n)$
- $\bullet \text{ space } S(n) = O(n)$

• ignoring prio points, we have a search tree!



• A treap can answer queries of the type $[x_{\min}, x_{\max}] \times [y_{\min}, \infty]$ • treap = tree and heap at the same time (41) (22) (39) (17) (13) (20)

Construction:

- \bullet extract max-prio point and x-median point
- recursively construct left and right subtrees
- $\bullet \text{ time } P(n) = O(n \log n)$
- space S(n) = O(n)Query:
- . . . the obvious way



• A treap can answer queries of the type $[x_{\min}, x_{\max}] \times [y_{\min}, \infty]$ • treap = tree and heap at the same time (41) (22) (39) (17) (13) (20)

Construction:

- \bullet extract max-prio point and x-median point
- recursively construct left and right subtrees
- $\bullet \text{ time } P(n) = O(n \log n)$
- space S(n) = O(n)Query:
- . . . the obvious way

Use treap as secondary structure in a 1-dimensional range tree:

•
$$S(n) = O(n \log n)$$

•
$$Q(n) = O(log^2n + k)$$



Things to consider for 'Practical' Solutions on large Data Sets

• use only space *linear* in the number of data items

Example:

- $n = 2^{29}$ (around 500 million)
- $n \log n = 29 \cdot n$, that is, we need 29 times more space than just to store the data!

Things to consider for 'Practical' Solutions on large Data Sets

- use only space *linear* in the number of data items
- try to avoid explicit storage of connectivity information



Example:

- \bullet data item has 10 bytes (lat/lon/prio)
- every node in the tree must have incoming pointer (8 Bytes)
- connectivity information almost *doubles* space
- Better: implicit connectivity e.g. via array representation:
 - leftchild(i)= $2 \cdot i + 1$
 - rightchild(i)= $2 \cdot i + 2$
 - parent(i)=i/2

Challenge:

- need 'almost complete' search tree
- requires more deliberate splitting choice than median

Things to consider for 'Practical' Solutions on large Data Sets

- use only space *linear* in the number of data items
- try to avoid explicit storage of connectivity information
- try to keep things in memory consecutively

More concrete:

- storing and scanning elements linearly in an array is surprisingly fast (cache effects)
- tree search issues rather dispersed memory accesses
- avoid many individual 'objects' allocated on the heap
- often a combination of a 'real' data structure with linear search very effective

Linear Scan:

- store (lat,lon,prio) consecutively in an array
- $Q(n) = \Theta(n)$ (always!)

Linear Scan:

- store (lat,lon,prio) consecutively in an array
- $Q(n) = \Theta(n)$ (always!)

Prioritized Linear Scan:

- store (lat,lon,prio) consecutively in an array **sorted** according to prio
- $Q(n) = \Theta(n')$ where n' is # of elements with high enough priority (overall)
- not bad if querying mainly large regions
- inefficient for very small regions

Linear Scan:

- store (lat,lon,prio) consecutively in an array
- $Q(n) = \Theta(n)$ (always!)

Prioritized Linear Scan:

- store (lat,lon,prio) consecutively in an array **sorted** according to prio
- $Q(n) = \Theta(n')$ where n' is # of elements with high enough priority (overall)
- not bad if querying mainly large regions
- inefficient for very small regions

Grid + Prioritized Linear Scans

- create grid
- store grid cells consecutively in memory
- within each grid cell store items sorted according to prio
- for query:
 - determine grid cells to be inspected
 - perform prioritized linear scan in each of them
- grid: limits effort if query region small
- prio sort: limits effort if min priority high
- parameter to choose: cell size
- good, if queries to be expected are somewhat uniform, otherwise kd-prio-tree might pay off



 $b^5 \ a^4 \ d^6 \ c^4 \ f^9 \ e^1 \ h^9 \ i^3 \ g^0 \ k^4 \ j^1 \ l^7 \ m^5$

Linear Scan:

- store (lat,lon,prio) consecutively in an array
- $Q(n) = \Theta(n)$ (always!)

Prioritized Linear Scan:

- store (lat,lon,prio) consecutively in an array **sorted** according to prio
- $Q(n) = \Theta(n')$ where n' is # of elements with high enough priority (overall)
- not bad if querying mainly large regions
- inefficient for very small regions

Grid + Prioritized Linear Scans

- create grid
- store grid cells consecutively in memory
- within each grid cell store items sorted according to prio
- for query:
 - determine grid cells to be inspected
 - perform prioritized linear scan in each of them
- grid: limits effort if query region small
- prio sort: limits effort if min priority high
- parameter to choose: cell size
- good, if queries to be expected are somewhat uniform, otherwise kd-prio-tree might pay off

You will implement this in the exercises!



 $b^5 \ a^4 \ d^6 \ c^4 \ f^9 \ e^1 \ h^9 \ i^3 \ g^0 \ k^4 \ j^1 \ l^7 \ m^5$

Thank you

... for your attention!

Questions?

Applications? (\rightarrow also talk to Filip!) See you in the exercise session ...

How bad can kd-Tree query time become?

How bad can kd-Tree query time become?

... in fact pretty bad:



In this perturbed 8×8 grid we have n = 64 and the query box intersects $\Omega(\sqrt{n})$ many cells without reporting anything.

$$\Rightarrow S(n) = \Omega(\sqrt{n} + k)$$

Can it become worse?

Can it become worse?

Fortunately, NO! We can prove a bound of

 $Q(n) = O(\sqrt{n} + k)$

due to the following observations:

- rectangles fully contained in query are already accounted for in the output size k
- only need to count rectangles properly *intersected* by query rectangle
- even simpler: consider the number of rectangles intersected by a vertical line

Can it become worse?

Fortunately, NO! We can prove a bound of

 $Q(n) = O(\sqrt{n} + k)$

due to the following observations:

- rectangles fully contained in query are already accounted for in the output size k
- only need to count rectangles properly *intersected* by query rectangle
- even simpler: consider the number of rectangles intersected by a vertical line

How many cells can be intersected by one vertical line?

- assume all nodes on even levels correspond to an X-split, on odd levels to an Y-split; root has level 0
- let T(n) be the number of cells intersected below an X-split node with n nodes in its subtree



The following holds:

- $\bullet \ T(n) = 1 + 1 + 2 \cdot T(n/4)$
- $\bullet T(1) = 1$

Can it become worse?

Fortunately, NO! We can prove a bound of

 $Q(n) = O(\sqrt{n} + k)$

due to the following observations:

- rectangles fully contained in query are already accounted for in the output size k
- only need to count rectangles properly *intersected* by query rectangle
- even simpler: consider the number of rectangles intersected by a vertical line

How many cells can be intersected by one vertical line?

- assume all nodes on even levels correspond to an X-split, on odd levels to an Y-split; root has level 0
- let T(n) be the number of cells intersected below an X-split node with n nodes in its subtree



The following holds:

- $\bullet \ T(n) = 1 + 1 + 2 \cdot T(n/4)$
- T(1) = 1

This yields the following sum with $(\log n)/2$ summands:

$$2 + 4 + 8 + \dots = \sum_{i=1}^{(\log n)/2} 2^{i}$$
$$< 2^{(\log n)/2+1} = 2 \cdot \sqrt{n}$$

Can it become worse?

Fortunately, NO! We can prove a bound of

 $Q(n) = O(\sqrt{n} + k)$

due to the following observations:

- rectangles fully contained in query are already accounted for in the output size k
- only need to count rectangles properly *intersected* by query rectangle
- even simpler: consider the number of rectangles intersected by a vertical line

How many cells can be intersected by one vertical line?

- assume all nodes on even levels correspond to an X-split, on odd levels to an Y-split; root has level 0
- let T(n) be the number of cells intersected below an X-split node with n nodes in its subtree



The following holds:

- $\bullet \ T(n) = 1 + 1 + 2 \cdot T(n/4)$
- $\bullet T(1) = 1$

This yields the following sum with $(\log n)/2$ summands:

$$2 + 4 + 8 + \dots = \sum_{i=1}^{(\log n)/2} 2^{i}$$
$$< 2^{(\log n)/2+1} = 2 \cdot \sqrt{n}$$

$$\Rightarrow Q(n) = O(\sqrt{n} + k)$$