

Some useful snippets!

The Auton Lab
Carnegie Mellon University



www.autonlab.org

Priebe/BIRCH/Theisson

- Suppose you want to do clustering but can only afford one pass through the data.
- Step 1: Choose an epsilon: the radius of your cluster.
- Step 2: Slog through the data.

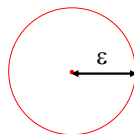
Priebe/BIRCH/Theisson



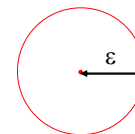
Priebe/BIRCH/Theisson

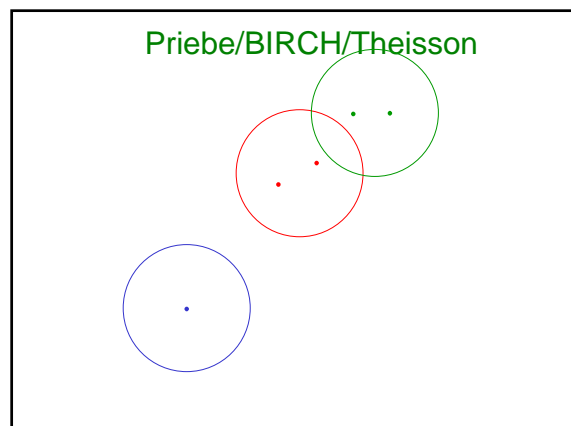
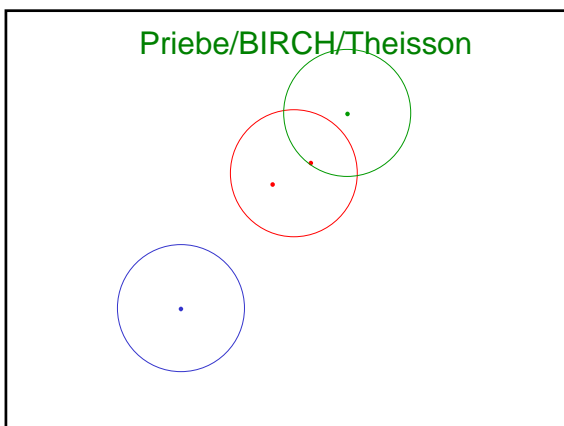
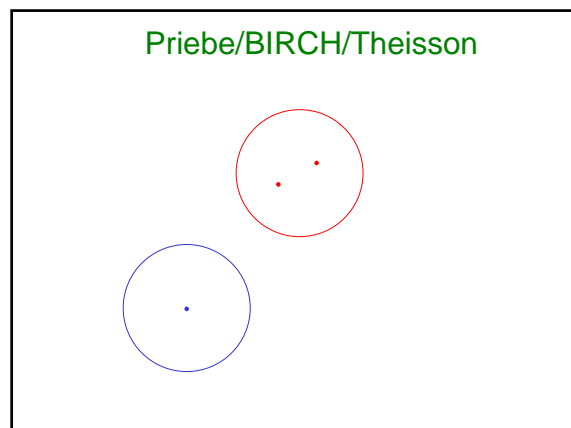
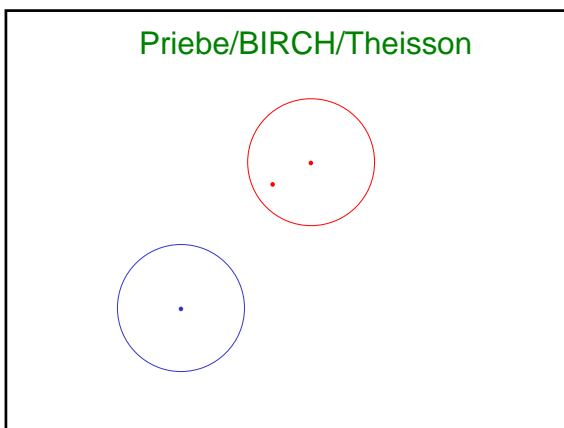
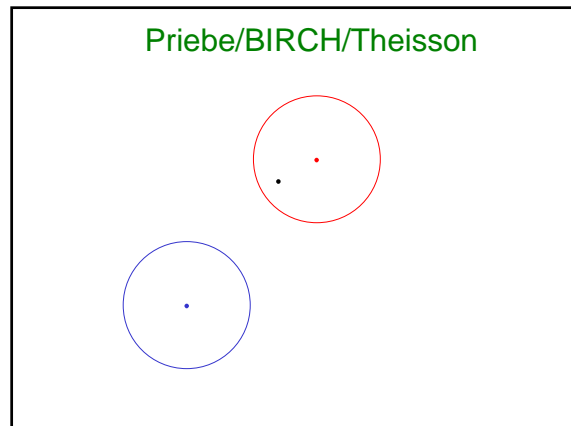
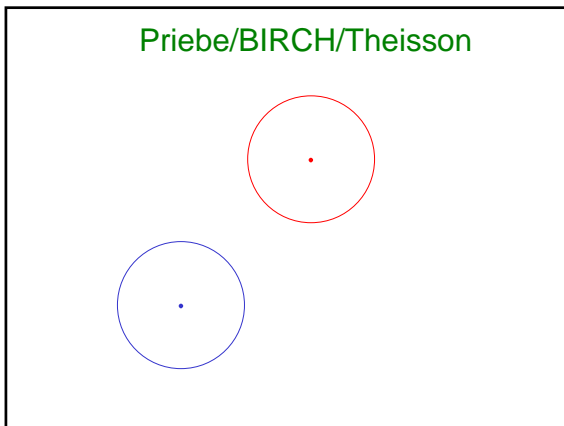


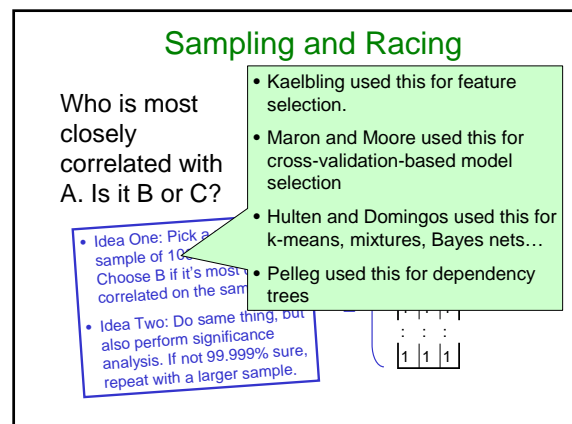
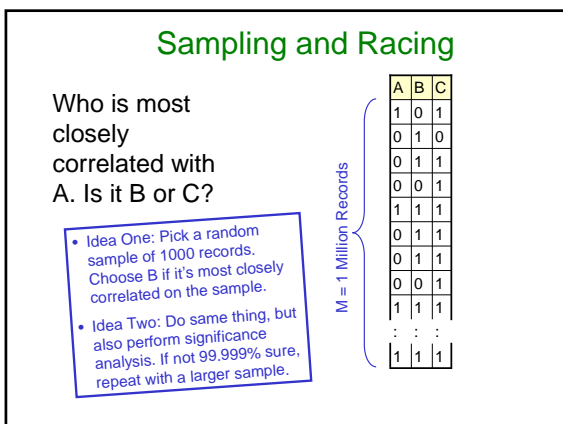
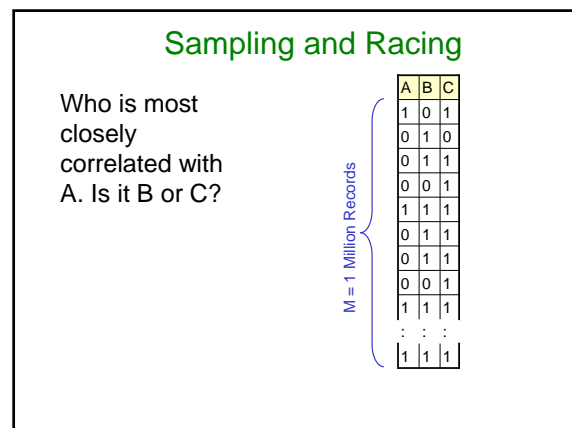
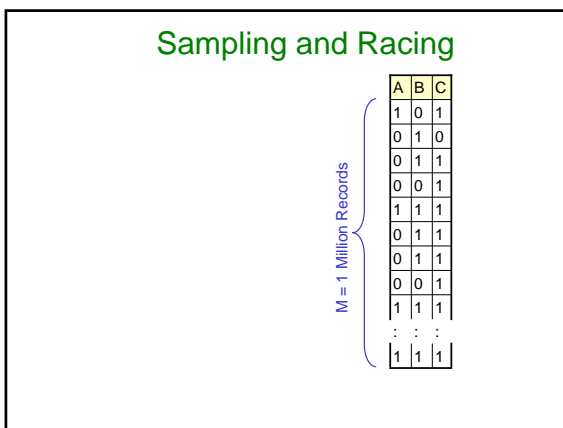
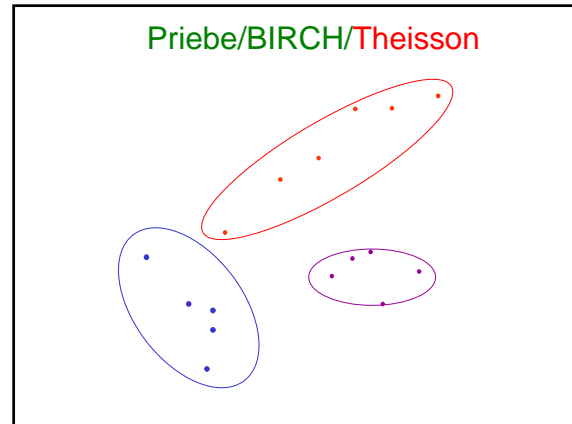
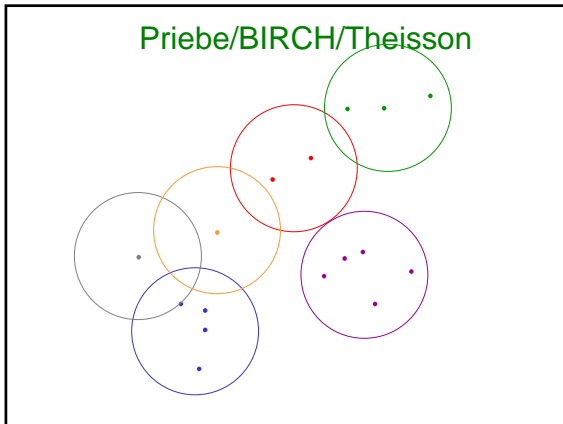
Priebe/BIRCH/Theisson



Priebe/BIRCH/Theisson







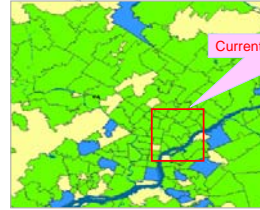
One Step of Spatial Scan

Entire area being scanned



One Step of Spatial Scan

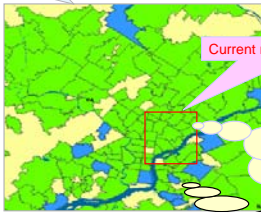
Entire area being scanned



Current region being considered

One Step of Spatial Scan

Entire area being scanned



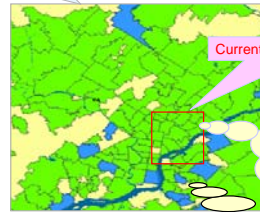
Current region being considered

I have a population of 5300 of whom 53 are sick (1%)

Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)

One Step of Spatial Scan

Entire area being scanned



Current region being considered

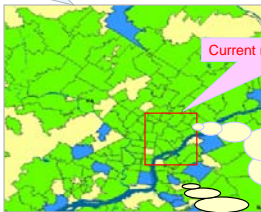
I have a population of 5300 of whom 53 are sick (1%)

Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)

So... is that a big deal?
Evaluated with Score function (e.g. Kulldorf's score)

One Step of Spatial Scan

Entire area being scanned



Current region being considered

I have a population of 5300 of whom 53 are sick (1%)

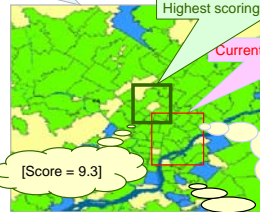
[Score = 1.4]

Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)

So... is that a big deal?
Evaluated with Score function (e.g. Kulldorf's score)

Many Steps of Spatial Scan

Entire area being scanned



Highest scoring region in search so far

Current region being considered

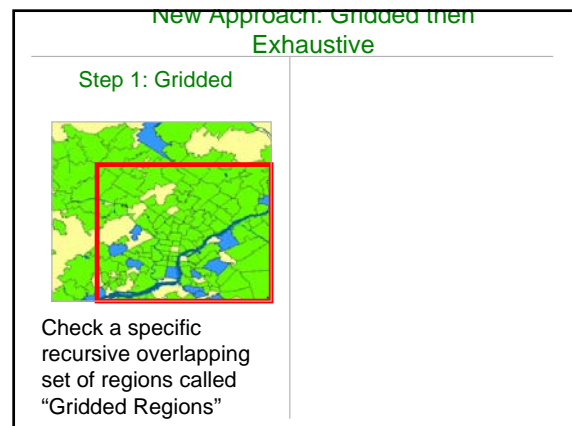
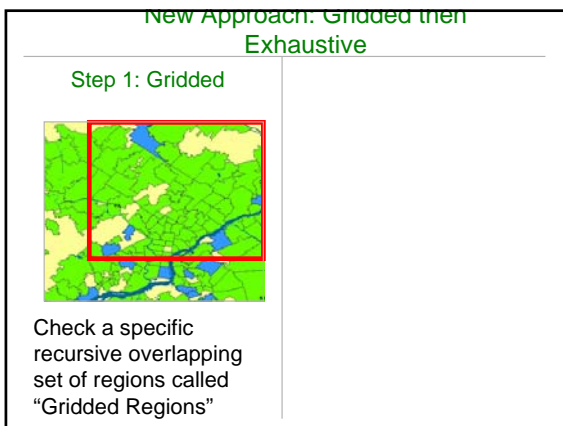
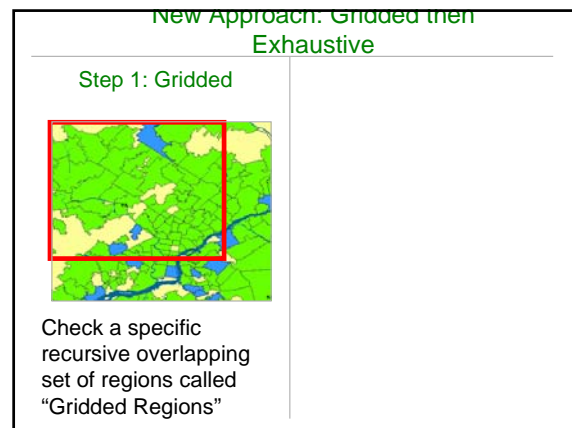
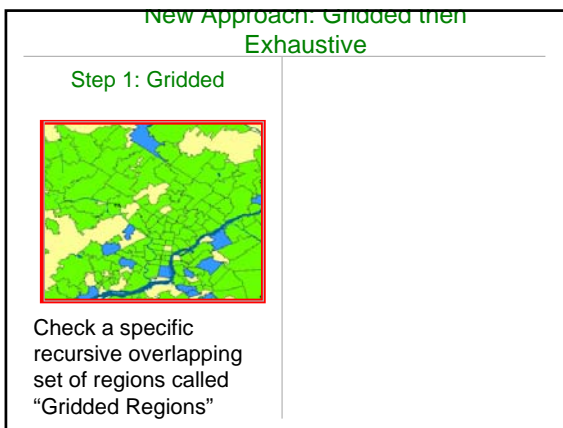
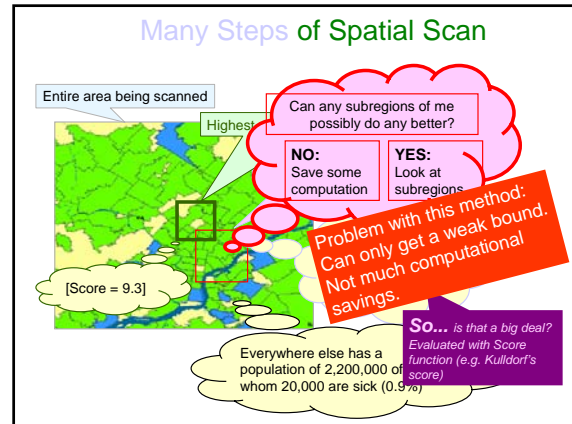
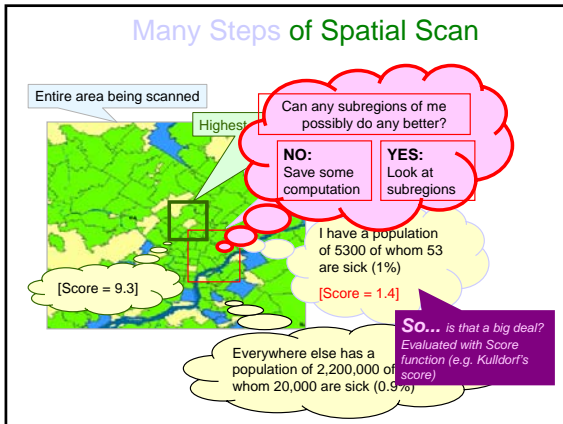
I have a population of 5300 of whom 53 are sick (1%)

[Score = 1.4]

Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)


So... is that a big deal?
Evaluated with Score function (e.g. Kulldorf's score)

[Score = 9.3]



New Approach: Gridded then Exhaustive


Step 1: Gridded



Check a specific recursive overlapping set of regions called "Gridded Regions"

New Approach: Gridded then Exhaustive


Step 1: Gridded



Check a specific recursive overlapping set of regions called "Gridded Regions"

New Approach: Gridded then Exhaustive


Step 1: Gridded



Check a specific recursive overlapping set of regions called "Gridded Regions"

New Approach: Gridded then Exhaustive


Step 1: Gridded



Check a specific recursive overlapping set of regions called "Gridded Regions"

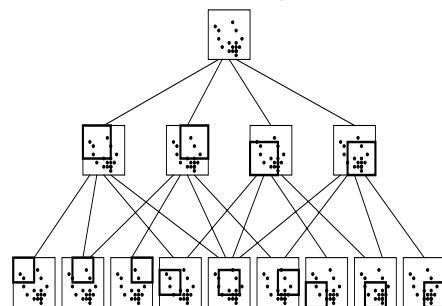
New Approach: Gridded then Exhaustive

Step 1: Gridded




Check a specific recursive overlapping set of regions called "Gridded Regions"

The multi-resolution tree for square regions



New Approach: Gridded then Exhaustive


Step 1: Gridded



Check a specific recursive overlapping set of regions called "Gridded Regions"

Step 2: Exhaustive

Consider the set of subregions of a Gridded Region.




New Approach: Gridded then Exhaustive

A subregion of me could be one of five types...

- ...entirely inside my top left gridded child
- ...entirely inside my top right gridded child
- ...entirely inside my bottom left gridded child
- ...entirely in my bottom right gridded child
- ...not entirely inside any of my 4 gridded children

Step 2: Exhaustive

Consider the set of subregions of a Gridded Region.



New Approach: Gridded then Exhaustive

A subregion of me could be one of five types...

- ...entirely inside my top left gridded child
- ...entirely inside my top right gridded child
- ...entirely inside my bottom left gridded child
- ...entirely in my bottom right gridded child
- ...not entirely inside any of my 4 gridded children

Step 2: Exhaustive

Consider the set of subregions of a Gridded Region.

FACT: Any subregion of this type must be big...

...and we can put fairly tight bounds on how well any region of this type can score

New Approach: Gridded then Exhaustive

A subregion of me could be one of five types...

- ...entirely inside my top left gridded child
- ...entirely inside my top right gridded child
- ...entirely inside my bottom left gridded child
- ...entirely in my bottom right gridded child
- ...not entirely inside any of my 4 gridded children

Procedure: Exhaust(Gridded Region)

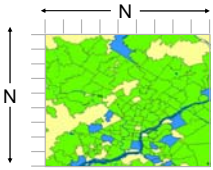
1. Exhaust(Region.TopLeft)
2. Exhaust(Region.TopRight)
3. Exhaust(Region.BottomLeft)
4. Exhaust(Region.BottomRight)
5. Is it possible that any "Type 5" subregion of "Gridded Region" could score better than best known score to date?

NO: Quit Procedure!

YES: Check all "Type 5" Subregions

region of this type can score

Fast squares speedup



- Theoretical complexity of fast squares: $O(N^2)$ (as opposed to naïve N^3), if maximum density region sufficiently dense.
- If not, we can use several other speedup tricks.
- In practice: 10-200x speedups on real and artificially generated datasets.

Emergency Dept. dataset (600K records): 20 minutes, versus 66 hours with naïve approach.