MapReduce Based Spatial Hotspots Detection Using Polygon Propagation

Jian Chen ^{1,2}, Satya Katragadda ², Shaaban Abbady ²

1: Department of Geography, University of North Alabama 2: NSF Center for Visual & Decision Informatics (CVDI), University of Louisiana at Lafayette

2019 AAG Annual Meeting, Washington, D.C.

Symposium on Frontiers in Geospatial Data Sciences: Big Data Computing for Geospatial Applications April 4th, 2019





Outlines

- Background
 - NSF Center fro Visual & Decision Informatics (CVDI)
 - MapReduce
 - Hotspot Analysis
- Our Approach
 - Algorithms
 - Experiments
- Conclusion





Background











MapReduce

- A processing technique and a program model for distributed computing based on Java.
- Allows massive scalability cross hundreds or thousands of servers in a Hadoop cluster.
 MAP REDUCE
- Mappers

Center for Visual & Decision Informatics





Why hotspots?

Spatial Hotspots

 Areas of space with unusually high incidence of events



Crash hotspots



Terrorism hotspots





Seismic hotspots





Center for Visual & Decision Informatics

Hotspot Analysis

- Hotspot analysis spatial clustering/spatial autocorrelation
- Research question How to know where interesting events (hotspots) occur with overwhelming data compilation?



- Existing approaches
 - Moran's I (Moran 1948), Getis-Ord general G (Getis & Ord 1992)
 - Getis-Ord Gi* (Getis 2007), Anselin's LISA (Anselin 1995)
 - Scan Statistics (Kulldorff 1997, 2006)







State-of-the-art: Scan Statistics

- Identifies clusters of high activity using spatiotemporal data
- Good for discrete localized outbreaks
- Can be combined with others for multivariate analysis
- Requires pre-defined boundaries
- Resultant hotspots w/ artificial shape may include non-hotspots (high false positive errors)



Hotspot analysis in public health, source: ESRI



Terrorism activities hotspots in Philippines Mar – Jun 2002, sources (Gao, et.al, 2013)



Center for Visual & Decision Informatics

Our Approach





- Aim:
 - Unsupervised
 - Absence of geographical boundary
 - Identifies compact hotspots
- Approach:
 - Density-based
 - Polygon propagation
- Details: Katragadda et. al, 2018





- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot





Given a set of points

- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon

Center for Vis

- Calculate log-Likelihood value
- If log-likelihood value > t2, it is categorized as hotspot





- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot







- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot







- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot







- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot







- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot







- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot







- Identify a triangle: density > t
 - Add nearest point to polygon: density > t
 - Until no more points can be added
- Until no more triangles can be formed
- For each Polygon
 - Calculate log-Likelihood value
 - If log-likelihood value > t2, it is categorized as hotspot









Center for Visual & Decision Informatics



Challenge of Polygon Propagation based Methods

- Creating a patch of areas covering the whole map
 - ✓ Using scan statistics to identify most likely cluster
- Polygons may spread to a large area W/O constraints
 - \checkmark A compactness parameter introduced to limit the size of the cluster P – the polygon, E- the smallest enclosing eclipse $\checkmark \alpha = \frac{Area(P)}{Area(E)}$
 - compasses all the points in polygon P

✓ Prioritize polygon propagate within the encompassing eclipse

- Polygon propagation is expensive, complexity $-O(N^3 log N)$
 - ✓ All possible polygons are pre-computed by Delaunay Triangulation (O(NlogN))
 - ✓ Distributed polygon propagation





Our Hotspots Detection Algorithm – MapReduce



Center for Visual & Decision Informatics



Given a set of input points

Mapper Create/extend polygons: density of polygon>t1 Reducer re-compute the polygons

Repeat until stabilization





Given a set of input points

Mapper Create/extend polygons: density of polygon>t1 Reducer re-compute the polygons

Repeat until stabilization

Center for Visu





Given a set of input points

Mapper Create/extend polygons: density of polygon>t1 Reducer re-compute the polygons

Repeat until stabilization







Given a set of input points

Mapper Create/extend polygons: density of polygon>t1 Reducer re-compute the polygons

Repeat until stabilization







Given a set of input points

Mapper Create/extend polygons: density of polygon>t1 Reducer re-compute the polygons

Repeat until stabilization





Given a set of input points

Mapper Create/extend polygons: density of polygon>t1 Reducer re-compute the polygons

Repeat until stabilization





Run Time Comparison

- Serial Implementation
 - CPU: Quad-Core 2.5 GHz
 - 8 GB RAM

- Hadoop Implementation (3-node cluster of)
 - CPU: Dual Six-Core 2.0
 GHz
 - 16 GB RAM

9 hours 32 minutes

– 48 minutes

Run time reduced about 90%!





Hotspots Detection on a Simulated Dataset



Original data





Elliptical scan statistics



Circular scan statistics







Hotspots Detection on Smiley Dataset



Original data

Circular scan statistics

Our method

Significant reduction in false positives – a major advantage





Test on Public Health Dataset

New York Breast Cancer Data

- Sources: NYSDH 2015, Bosceo et. al, 2016
- Breast cancer cases aggregated to census block
- 13,848 points (centroid of a census block), representing 72, 926 patients and a population of 27, 820, 632 in 2009 in the State of New York







Test on Public Health Dataset





Conclusion

- Polygon propagation (PP) can detect better heterogeneous cluster that are irregularly shaped
- PP detects more compact hotspots area (has less false positive errors) than the state-of-the-art Scan Statistics
- MapReduce based PP significantly outperforms the traditional Scan Statistics and PP algorithms, make it suitable for largescale spatial hotspots detection
- Future work
 - Identify better ways to calculate the stability of the detected hotspot
 - Reduce the computational time of the greedy search



Acknowledgement

• NSF Grants CNS-16505511 and IIP-1160958

• NSF I/UCRC CVDI IAB grant CS-14-4

University Research Grant, University of North Alabama











References

- Anselin, L. 1995. "Local Indicators of Spatial Association-LISA." Geographical Analysis 27 (2): 93–115.
- Boscoe, F., T. Talbo, and M. Kulldorff. 2016. "Public Domain Small-area Cancer Incidence Data for New York State, 2005–2009." Geospatial Health 11 (1): 304.
- Getis, A., and J. K. Ord. 1992. "The Analysis of Spatial Association by Use of Distance Statistics." Geographical Analysis 24 (3): 189–206.
- Getis, A. 2007. "Reflections on Spatial Autocorrelation." Regional Science and Urban Economics 37: 491–496.
- Gao, P, P. Guo, K. Liao, J. J. Webb, and S. L. Cutter. 2013. "Early detection of terrorism outbreaks using prospective space-time scan statistics." *The Professional Geographer* 65(4): 676-691.
- Katragadda, S., J. Chen, and S. Abbady. 2018. "Spatial Hotspot Detection using Polygon Propagation", International Jornal of Digital Earth, DOI: 10.1080/17538947.2018.1485754
- Kulldorff, M. 1997. "A Spatial Scan Statistic." Communications in Statistics Theory and Methods 26 (6): 1481–1496.
- Kulldorff, M., L. Huang, L. Pickle, and L. Duczmal. 2006. "An Elliptic Spatial Scan Statistic." Statistics in Medicine 25 (22): 3929–3943.
- Moran, P. A. P. 1950. "Notes on Continuous Stochastic Phenomena." Biometrika 37 (1/2): 17–23.
- NYSDH (New York State Department of Health). 2015. "Cancer Mapping Data: 2005-2009", <u>https://health.data.ny.gov/Health/Cancer-Mapping-Data-2005-2009/cw3n-</u> <u>fkji?category=Health&view_name=Cancer-Mapping-Data-2005–2009</u>
- Ord, J. K., and A. Getis. 2001. "Testing for Local Spatial Autocorrelation in the Presence of Global Autocorrelation." Journal of Regional Science 41 (3): 411–432.



Center for Visual & Decision Informatics

Questions & Comments





