Faculty Development Programme on "Applications of Artificial Intelligence on Geospatial Data"
sponsored by AICTE Training and Learning Academy, July 26th-30th,, 2021
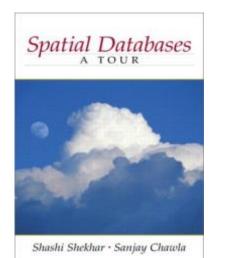Maulana Abul Kalam Azad University of Technology, and State Technological University of West Bengal, India
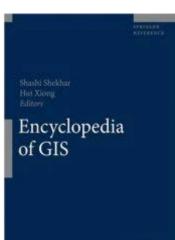
# What is Special about Spatial Data Science and GeoAI?

## Shashi Shekhar

McKnight Distinguished University Professor, University of Minnesota

www.cs.umn.edu/~shekhar, shekhar@umn.edu

Spatial Databases A TOUR — Shashi Shekhar · Sanjay Chawla



Shashi Shekhar, Hui Xiong, Editors — Encyclopedia of GIS — Springer



University Consortium for GEOGRAPHIC INFORMATION SCIENCE — A UCGIS Call to Action: Bringing the Geospatial Perspective to Data Science Degrees and Curricula — Summer 2018



SPATIAL COMPUTING — SHASHI SHEKHAR AND PAMELA VOLD — THE MIT PRESS ESSENTIAL KNOWLEDGE SERIES

# A Spatial Data Science Story

<u>1854: What causes Cholera?</u>

Miasma theory

```
Collect &        →   Discover Patterns,   →   Test Hypothesis    →   Develop
Curate Data          Generate Hypothesis      (Experiments)          Theory
```

? water pump

Remove pump handle

TURNING POINTS IN SCIENCE
GERM THEORY

- Pump sites
- Deaths from cholera

The FOURTH PARADIGM
DATA-INTENSIVE SCIENTIFIC DISCOVERY
EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

**Impact:** hygiene, drinking water supply, sewage system, …

nature
BIG DATA
SCIENCE IN THE PETABYTE ERA

**Q? What are Choleras of today?**
**Q? How may Spatial Data Sc. Help?**

# What has changed? **Spatial Data Revolution**

| Spatial | Last Century | Last Decade |
|---|---|---|
| **Data** | Few satellites and sensors | Nano-satellites, Billions of GPS enabled smartphones |
| **Data Access** | Need special hardware and network | |
| **Spatial Platforms** | ESRI Arc/Info | |
| **Spatial Data Science** | Spatial Patterns, e.g., hotspots (SatScan, ESRI Geostatistics) | |
| **Spatial Visualization** | Quilt, e.g., MS Terraserver | |

# Spatial Computing is Ubiquitous Today!

- 2 billion GPS receivers in use, will hit 7 billion by 2022.

- Besides location, it reference time for critical infrastructure
    - Telecommunications industry
    - Banks
    - Airlines...

- GPS is the single point of failure for the entire modern economy.

- 50,000 incidents of deliberate (GPS) jamming last two years
    - Against Ubers, Waymo's self-driving cars, delivery drones from Amazon

**Bloomberg Businessweek**

July 25, 2018, 4:00 AM CDT

The World Economy Runs on GPS. It Needs a Backup Plan

*Source:* *https://www.bloomberg.com/news/features/2018-07-25/the-world-economy-runs-on-gps-it-needs-a-backup-plan*

# Large Constellations of Small Satellites

- Hi-frequency (e.g., daily or hourly) time-series of imagery of entire earth
- Small Satellites: video (5-minutes): https://geospatialstream.com/sciencecasts-nasa-embraces-small-satellites/
- Large Constellations
  - 2021: Planet Labs: 200+ satellites: daily Earth scan (1m resolution, visible+NIR bands)

Source: WorldView FAQ, blog.digitalglobe.com/news/frequently-asked-questions-about-worldview-4/



| 540 cm/212.4 in | | | | | | |
| 335 cm/131.89 in | | | | | | |
| 177 cm/69.7 in | | | | | | |
| 117 cm/46.06 in | | | | | | |
| 80 cm/31.5 in | | | | | | |
| 75 cm/29.53 in | | | | | | |
| 30 cm/11.88 in | | | | | | |
| Human | Planet Labs | BlackSky | Terra Bella | BlackBridge | Pleiades 1B | DigitalGlobe WorldView-4 |

Spatial Computing
Research Group

# Spatial Data Revolution

1. **GPS & Location traces**
   - 2 billion GPS receivers today (7 billion by 2022)
   - Reference clock for telecom, banks, …
   - Help understand Spatio-temporal patterns of life

2. **(Nano-)Satellite Imagery**, …



The World Economy Runs on GPS. It Needs a Backup Plan

**Bloomberg Businessweek**

July 25, 2018, 4:00 AM CDT



**ENSURING RESOURCE AVAILABILITY**

Advanced technology, including many types of Earth information, will unlock up to *$1.6 trillion* in economic savings for energy generation and use by 2035.

Satellite observations can also help ensure water availability, which is particularly important to the 20% of the world now living in areas of water scarcity.
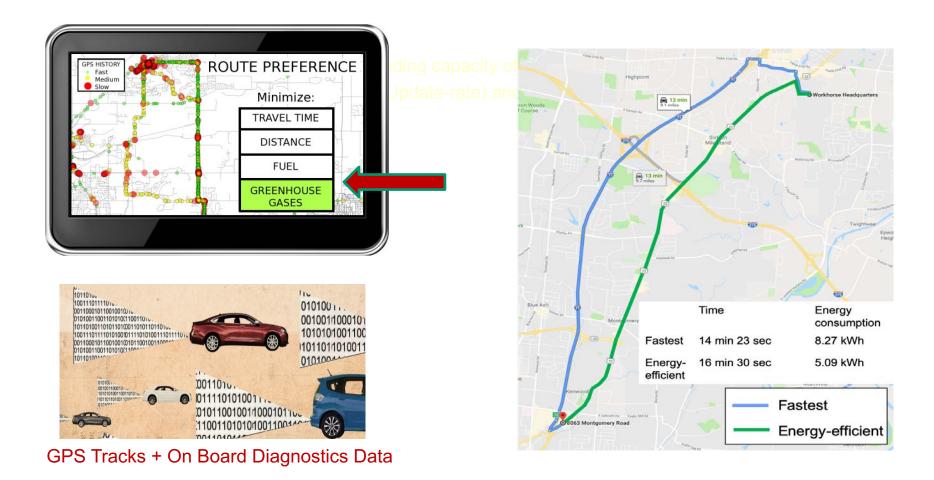


McKinsey Global Institute

The study estimates that the use of personal location data could save consumers worldwide more than $600 billion annually by 2020. Computers determine users' whereabouts by tracking their mobile devices, like cellphones.

The New York Times

Published: May 13, 2011

**Source:** Y. Xie et al., Transforming Smart Cities With Spatial Computing, Proc. IEEE Intl. Conf. on Smart Cities, 2018.

# Next Generation Navigation App
# to Reduce Emission, Energy use



GPS Tracks + On Board Diagnostics Data

| | Time | Energy consumption |
|---|---|---|
| Fastest | 14 min 23 sec | 8.27 kWh |
| Energy-efficient | 16 min 30 sec | 5.09 kWh |

**Details:** Y. Li, P. Kotwal, P. Wang, Y. Xie, S. Shekhar, and W. Northrop, Physics-guided Energy-efficient Path Selection Using On-board Diagnostics Data, ACM/IMS Trans. Data Sc. 1(3):1-28, Article 22, Oct. 2020.

# What has changed? **Spatial Data Access**

| Spatial | Last Century | Last Decade |
|---|---|---|
| Data | Few satellites and sensors | Nano-satellites, Billions of GPS enabled smartphones |
| Data Access | Need special hardware and network | Cloud based repositories and analytics (e.g., DARPA GCA) |
| Spatial Platforms | ESRI Arc/Info | |
| Spatial Data Science | Spatial Patterns, e.g., hotspots (SatScan, ESRI Geostatistics) | |
| Spatial Visualization | Quilt, e.g., MS Terraserver | |

# Easier Access: Cheap (or free) Cloud Repositories

- 2008: USGS gave away 35-year LandSat satellite imagery archive
  - Analog of public availability of GPS signal in late 1980s
- 2017: Cloud-based repositories of geospatial data
  - Explosion in machine learning on satellite imagery to map crops, water, buildings, roads, …

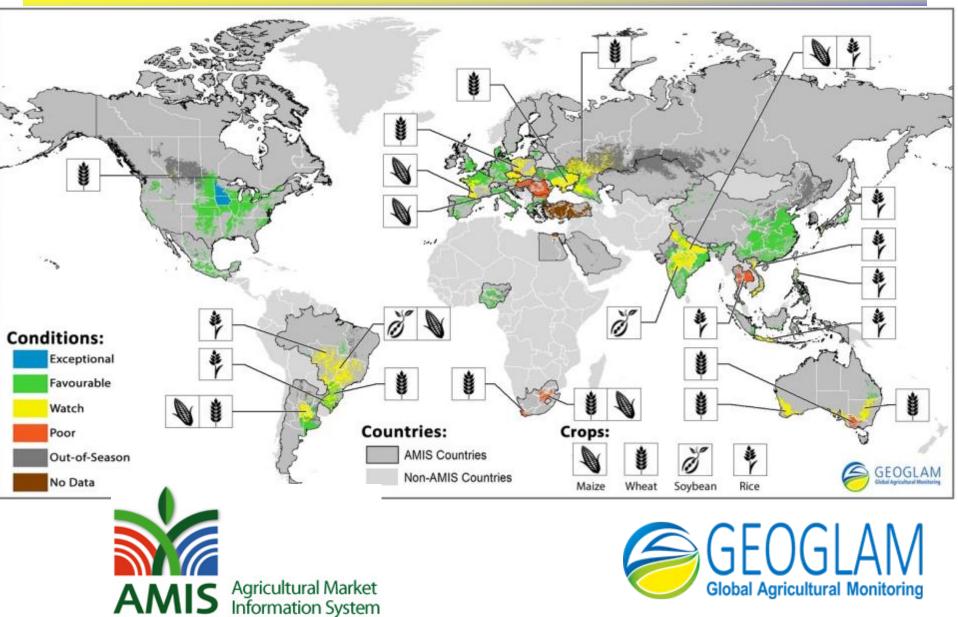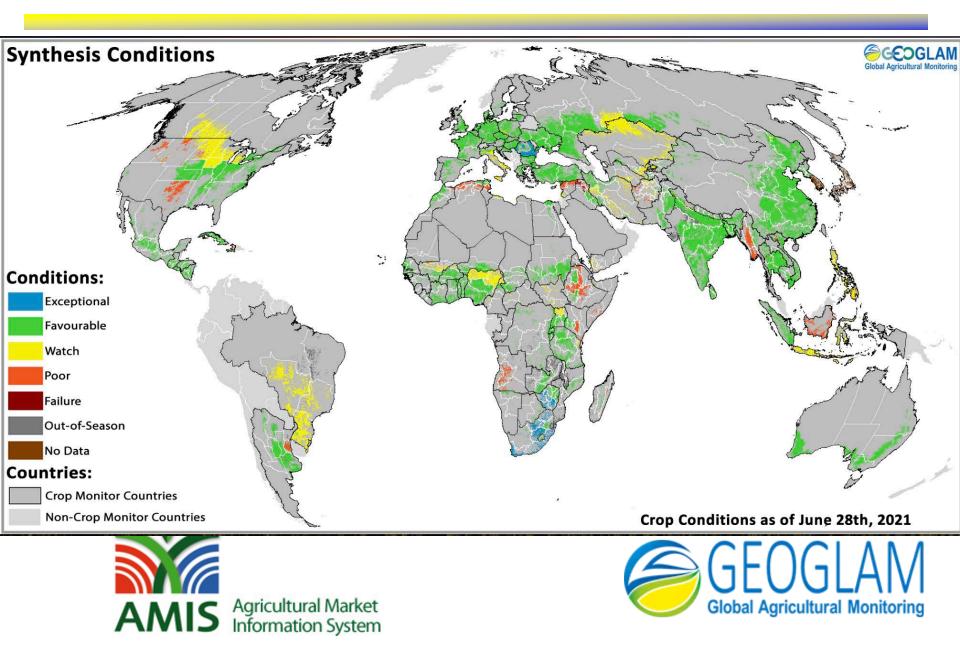|  | Google Earth Engines | NEX | AWS Earth |
|---|---|---|---|
| Elevation, Landsat, LOCA, MODIS, NAIP | x | x | x |
| NOAA | x |  | x |
| AVHRR, FIA, GIMMM, GlobCover, NARR, TRIMM, Sentinel-1 | x | x |  |
| IARPA, GDELT, MOGREPS, OpenStreetMap, Sentinel-2, SpaceNet (building/road labels for ML) |  |  | x |
| CHIRPS, GeoScience Australia, GSMap, NASS, Oxford Map, PSDI, WHRC, WorldClim, WorldPop, WWF, | x |  |  |
| BCCA, FLUXNET |  | x |  |

Google Earth Engine

NEX NASA Earth Exchange

Earth on AWS
Build planetary-scale applications in the cloud with open geospatial data.

Spatial Computing Research Group

# Global Agriculture Monitoring

# Global Agriculture Monitoring

# What has changed? **Spatial Big Data Platforms**

| Spatial | Last Century | Last Decade |
|---|---|---|
| **Data** | Few satellites and sensors | Nano-satellites, Billions of GPS enabled smartphones |
| **Data Access** | Need special hardware and network | Cloud based repositories, e.g., Earth on AWS |
| **Spatial Platforms** | ESRI Arc/Info, SQL3/OGC, e.g., Postgis, | Geospatial Cloud Analytics (Monitor crops, fracking, illegal fishing), ESRI GIS Tools for Hadoop, … |
| **Sp... Sci...** | | |
| **Sp... Vis...** | | |

# Spatial Data Types >> Points

**Q?** What is distance between Washington D.C. and U.S.A.?
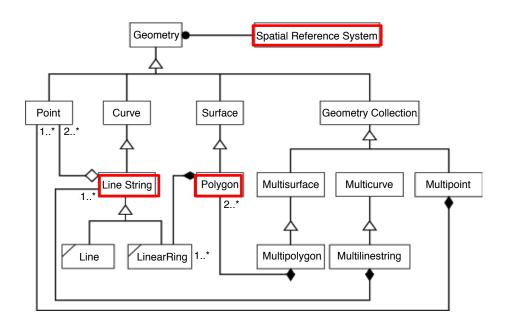
- Zero ( Washington D.C. is inside U.S.A. )
- NSF OKN funded 2 grants on geo-knowledge networks!

# Spatial Data Types: OGC Simple Features Standard

- Data types: Point, LineString, Polygon, Collections
- Relationships: Topological, Metric (e.g., distance)
- ML Challenge: implicit spatial relationships
  - Approach: pre-compute for feature extraction



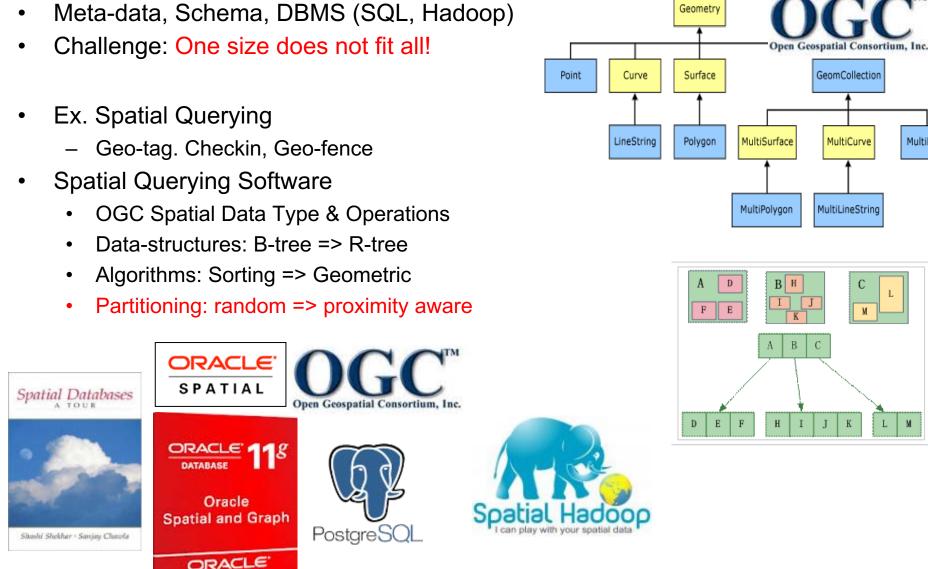| Basic Functions | SpatialReference () |
|---|---|
| | Envelop () |
| | Export () |
| | IsEmpty () |
| | IsSimple () |
| | Boundary () |
| Topological / Set Operators | Equal |
| | Disjoint |
| | Intersect |
| | Touch |
| | Cross |
| | Within |
| | Contains |
| | Overlap |
| Spatial Analysis | Distance |
| | Buffer |
| | ConvexHull |
| | Intersection |
| | Union |
| | Difference |
| | DymmDiff |

**Details:** S. Shekhar et al., Spatial Databases: Accomplishments and Research Needs, IEEE Trans. on Knowledge and Data Eng., 11(1):45-55, Jan.-Feb. 1999.

# Spatial Big Data Curation

- Meta-data, Schema, DBMS (SQL, Hadoop)
- Challenge: One size does not fit all!

- Ex. Spatial Querying
  - Geo-tag. Checkin, Geo-fence
- Spatial Querying Software
  - OGC Spatial Data Type & Operations
  - Data-structures: B-tree => R-tree
  - Algorithms: Sorting => Geometric
  - Partitioning: random => proximity aware

# Spatial Big Data Platforms

| Genre | Examples |
|---|---|
| Relational DBMS, Spatial Library | Oracle, IBM DB2, PostgreSQL, MS SQL Server OGC Simple Features, … |
| Parallel DBMS | Teradata, Vertica, Greenplum, DataAllegro, ParAccel |
| Big Data Platforms | Hadoop, MapReduce, Spark, Hbase, Hive, … |
| Spatial Big Data Platforms | ESRI GIS Tools for Hadoop, GeoWave, SpatialSpark, GeoSpark, Simba, Hadoop-GIS, SpatialHadoop, ST-Hadoop |

Blocks based on Hadoop HashPartition

Legend
- Block_0 (64M)
- Block_1 (64M)
- Block_2 (20M)

(a)

Blocks based on Spatial Data Partition

Legend
- Block_0 (38M)
- Block_1 (71M)
- Block_2 (39M)

(b)

**Source**: X.Yao et al., Computers and GeoScience, 106:60-67, 2017

# What has changed? **Spatial Data Science**

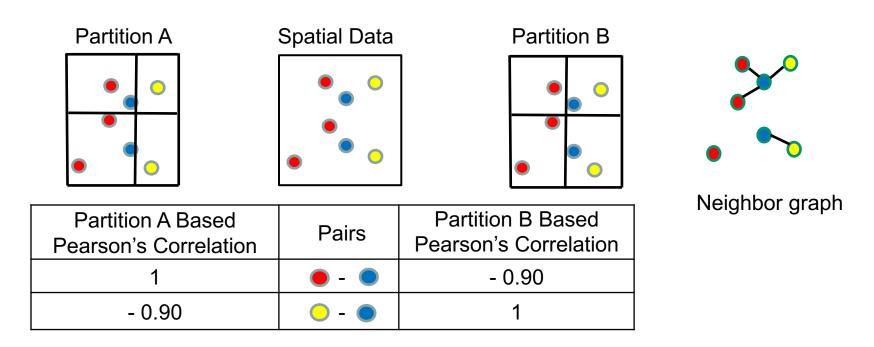| Spatial | Last Century | Last Decade |
|---|---|---|
| **Data** | Few satellites and sensors | Nano-satellites, Billions of GPS enabled smartphones |
| **Data Access** | Need special hardware and network | Cloud based repositories, e.g., Earth on AWS |
| **Spatial Platforms** | ESRI Arc/Info | SQL3/OGC, e.g., Postgis, ESRI GIS Tools for Hadoop, Google Earth Engine |
| **Spatial Data Science** | Spatial Patterns, e.g., hotspots (SatScan, ESRI Geostatistics) | (a) Spatial Network Patterns, e.g., linear hotspots<br>(b) Spatio-temporal (ST) patterns, e.g., Change time-series |
| **Spatial Visualization** | Quilt, e.g., MS Terraserver | |

# Spatial Challenges: Traditional Data Science

- Traditional methods not robust in face of
  - Challenge 1: Spatial continuity
  - Challenge 2: Auto-correlation, Heterogeneity , Edge-effect, …
  - Challenge 3: Noise

**Details:** [Data Science for Earth: The Earth Day Report](Data Science for Earth: The Earth Day Report) , E. Eftelioglu, S. Shekhar, J. Hudson, L. Joppa, C. Baru, V. Janeja, et al. ACM SIGKDD Explorations Newsletter, 22(1), May 2020

# Challenge 1: Continuous Space

- Traditional methods not robust in face of Spatial continuity
    - Gerrymandering risk: Classical methods not robust
    - Result changes if spatial partitioning changes
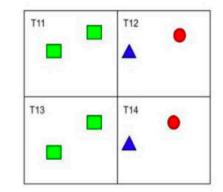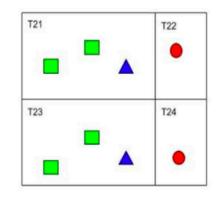    - Formally, Modifiable Areal Unit Problem (MAUP)



Partition A    Spatial Data    Partition B

Neighbor graph

| Partition A Based Pearson's Correlation | Pairs | Partition B Based Pearson's Correlation |
|---|---|---|
| 1 | 🔴 - 🔵 | - 0.90 |
| - 0.90 | 🟡 - 🔵 | 1 |

**Details:** [Data Science for Earth: The Earth Day Report](#) , E. Eftelioglu, S. Shekhar, J. Hudson, L. Joppa, C. Baru, V. Janeja, et al. ACM SIGKDD Explorations Newsletter, 22(1), May 2020
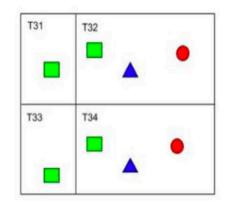
# Classical Data Mining Methods not robust either!

Consider the spatial Data in Figure (a)
Along with 3 alternative partitions in Figures (b), (c) and (d).
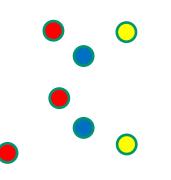


(a) Map of 3 item-types

(b) Spatial Partition P1

(c) Spatial Partition P2
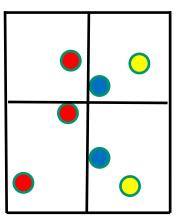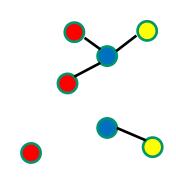
(d) Spatial Partition P3

| Spatial Partitioning Definition | P1 | P2 | P3 |
|---|---|---|---|
| Transactions | T11, T12, T13, T14 | T21, T22, T23, T24 | T31, T32, T33, T44 |
| **Associations** with support >= 0.5 | ( ▲ ● ) | ( ■ ▲ ) | ( ■ ▲ ● ) |

# Approach: Neighbor Graph

- Challenge: One size does not fit all

- Ex. Interaction patterns

(a) a map of 3 features          (b) Spatial Partitions          (c) Neighbor graph

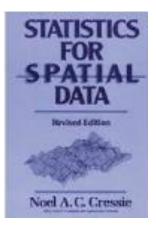|  | Pearson's Correlation | Ripley's cross-K | Participation Index |
|---|---|---|---|
| 🔴🔵 | -0.90 | 0.33 | 0.5 |
| 🟡🔵 | 1 | 0.5 | 1 |

Details:  Discovering Spatial Co-location Patterns: A General Approach,
IEEE Transactions on Knowledge and Data Eng., 16(12), December 2004 (w/ H.Yan, H.Xiong).

# A Metric of Spatial Cross-Correlation

- Ripley's Cross K-Function Definition

$$K_{ij}(h) = \lambda_j^{-1} E \text{ [number of type } \textit{\textbf{j}} \text{ event within distance } h \text{ of a randomly chosen type } \textit{\textbf{i}} \text{ event]}$$

  - Cross K-function of some pair of spatial feature types
  - Example
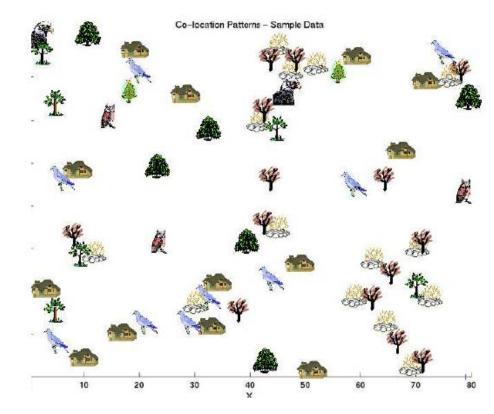    - Which pairs are frequently co-located
    - Statistical significance



STATISTICS FOR SPATIAL DATA

Revised Edition

Noel A. C. Cressie



Spatial Computing Research Group

# Co-locations

- Given: A collection of different types of spatial events
- Find: Co-located subsets of event types
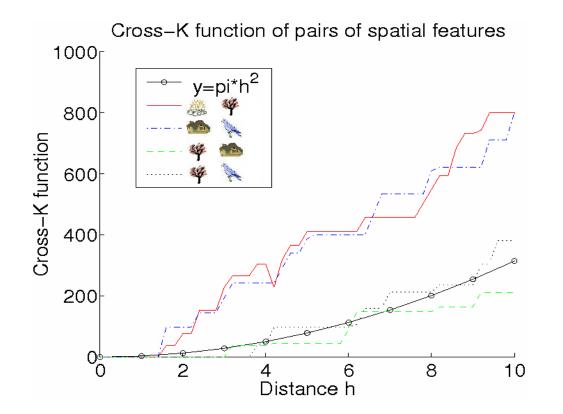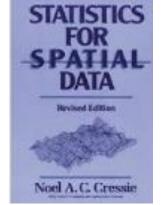




Co-location Patterns – Sample Data

Source: Discovering Spatial Co-location Patterns: A General Approach, IEEE Transactions on Knowledge and Data Eng., 16(12), December 2004 (w/ H.Yan, H.Xiong).

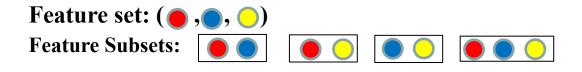Spatial Computing Research Group

# Illustration of Cross-Correlation

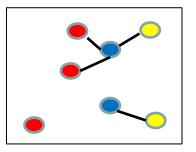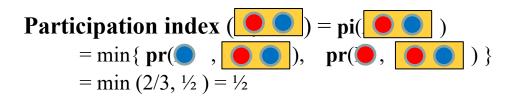- Illustration of Cross K-function for Example Data



Cross–K function of pairs of spatial features

# Spatial Colocation

**Feature set: (🔴 , 🔵 , 🟡 )**

**Feature Subsets:** [🔴 🔵]   [🔴 🟡]   [🔵 🟡]   [🔴 🔵 🟡]

**Participation ratio (pr):**

$\text{pr}($🔴 $,$ [🔴 🔵] $) =$ fraction of 🔴 instances neighboring feature {🔵} $= 2/3$

$\text{pr}($🔵 $,$ [🔴 🔵] $) = \frac{1}{2}$

**Participation index** ([🔴 🔵]) $= \text{pi}($[🔴 🔵]$)$

$= \min\{\ \text{pr}($🔵 $,$ [🔴 🔵]$),\quad \text{pr}($🔴 $,$ [🔴 🔵]$)\ \}$

$= \min(2/3, \frac{1}{2}) = \frac{1}{2}$

## Participation Index Properties:

(1) <u>Computational</u>: Non-monotonically decreasing like support measure

(2) <u>Statistical</u>: Upper bound on Ripley's Cross-K function

Source: Discovering Spatial Co-location Patterns: A General Approach, IEEE Transactions on Knowledge and Data Eng., 16(12), December 2004 (w/ H.Yan, H.Xiong).
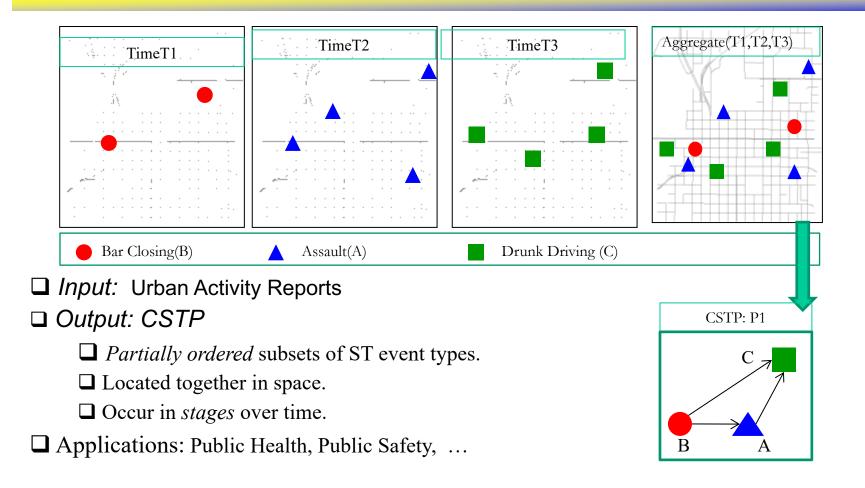
# Participation Index >= Cross-K Function



| | | | |
|---|---|---|---|
| **Cross-K (A,B)** | 2/6 = 0.33 | 3/6 = 0.5 | 6/6 = 1 |
| **PI (A,B)** | 2/3 = 0.66 | 1 | 1 |

# Spatial Colocation: Trends

- Algorithms
  - Join-based algorithms
    - One spatial join per candidate colocation
  - Join-less algorithms

- Statistical Significance
  - ?Chance-patterns

- Spatio-temporal
  - Which events co-occur in space and time?
    - (bar-closing, minor offenses, drunk-driving citations)
  - Which types of objects move together?

# Cascading spatio-temporal pattern (CSTP)



| ● | Bar Closing(B) | ▲ | Assault(A) | ■ | Drunk Driving (C) |

- ❑ *Input:* Urban Activity Reports
- ❑ *Output: CSTP*
    - ❑ *Partially ordered* subsets of ST event types.
    - ❑ Located together in space.
    - ❑ Occur in *stages* over time.
- ❑ Applications: Public Health, Public Safety, …
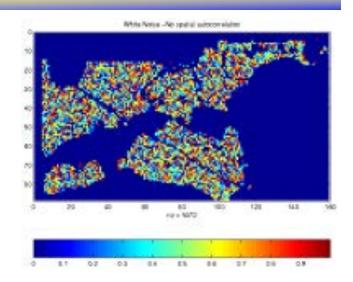
# Spatial Challenges: Traditional Data Science

- Traditional methods not robust in face of
  - Challenge 1: Spatial continuity
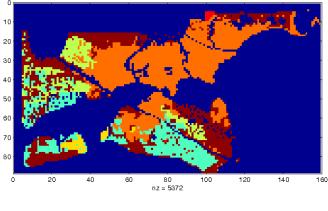  - Challenge 2: Auto-correlation, Heterogeneity , Edge-effect, …
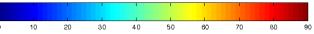  - Challenge 3: Noise

# Challenge 2: Spatial Auto-correlation

- Traditional Statistics, ML, Data Mining
- Ubiquitous i. i. d. assumption
    - Data samples independent of each other
    - From identical distribution
- Problem
    - Ignores auto-correlation, heterogeneity
    - Salt n Pepper noise

White Noise - No spatial autocorrelation

Vegetation distribution across the marshland

nz = 5372

STATISTICS FOR SPATIAL DATA

Revised Edition

Noel A. C. Cressie

Wiley Series in Probability and Statistics

Statistics for SPATIO-TEMPORAL DATA

Noel Cressie · Christopher K. Wikle

# Illustration of Location Prediction Problem



Nest Locations

Vegetation Index

Water Depth

Distance to Open Water

# Spatial Auto-Regression & Parameter Estimation

| Name | Model |
|------|-------|
| Classical Linear Regression | $\mathbf{y} = \mathbf{x}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ |
| Spatial Auto-Regression | $\mathbf{y} = \rho\mathbf{W}\mathbf{y} + \mathbf{x}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ |

$\rho$ : the spatial auto-regression (auto-correlation) parameter

$\mathbf{W}$ : $n$-by-$n$ neighborhood matrix over spatial framework

- **<u>Maximum Likelihood Estimation</u>**

$$\ln(L) = \boxed{\ln|\mathbf{I} - \rho\mathbf{W}|} - \frac{n\ln(2\pi)}{2} - \frac{n\ln(\sigma^2)}{2} - SSE$$

- Computing determinant of large matrix is a hard (open) problem!
    - size(W) is quadratic in number of locations/pixels.
    - Typical raster image has Millions of pixels
    - W is sparse but not banded.

**Details:** A parallel formulation of the spatial auto-regression model for mining large geo-spatial datasets, SIAM Intl. Workshop on High Perf. and Distr. Data Mining, 2004. (with B. Kazar)

# Comparing Traditional and Spatial Models

- Dataset: Bird Nest prediction
- Linear Regression
  - Lower prediction accuracy, coefficient of determination,
  - Residual error with spatial auto-correlation
- Spatial Auto-regression outperformed linear regression



ROC Curve for testing data(Stubble marshland 1995)

ROC Curve
for learning

Classical Regression
Spatial Regression



ROC Curve for learning data(Darr marshland 1995)

ROC Curve
for testing

Classical Regression
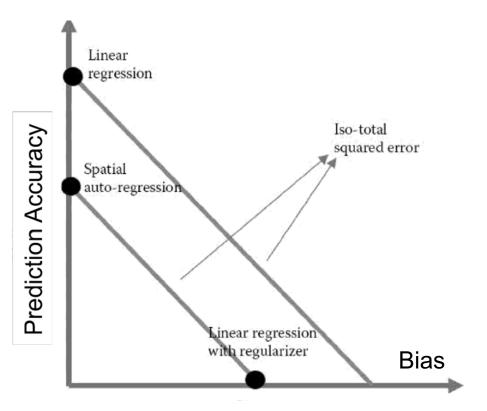Spatial Regression

Spatial Computing
Research Group

# Prediction Error and Bias Trade-off

- Linear Regression (LR)
$$y = X\beta + \varepsilon$$

- LR with Auto-correlation Regularizer
$$y = X\beta + \varepsilon$$
$$\varepsilon = \|y - X\beta\|^2 + \|y - y_{neighbor}\|^2$$

- Spatial Auto-Regression
$$y = \rho W y + X\beta + \varepsilon$$

# Research Needs for Location Prediction

- Spatial Auto-Regression
  - Estimate W
  - Scaling issue $\quad \rho \mathrm{W} y \text{ vs. } \mathrm{X}\beta$
- Spatial interest measure
  - e.g., distance(actual, predicted)

|  (a) | (b) | (c) | (d) | Legend |
|------|-----|-----|-----|--------|

(a) Actual Sites

(b) Pixels with actual sites

(c) Prediction 1

(d) Prediction 2.
Spatially more interesting than Prediction 1

Legend

$\odot$ = nest location

A = actual nest in pixel

P = predicted nest in pixel

Spatial Computing
Research Group

# Salt n Pepper Noise



wetland    dry land

Input:

Output:

train

test

(a) aerial photo    (b) aerial photo    (c) true classes    (d) DT prediction (Salt n Pepper Noise)    (e) SDT prediction

**Training samples**: upper half
**Test samples**: lower half
**Spatial neighborhood**: maximum 11 pixels by 11 pixels

DT: decision tree
SDT: spatial decision tree

Details: Focal-Test-Based Spatial Decision Tree Learning. IEEE Trans. Knowl. Data Eng. 27(6): 1547-1559, 2015 (summary in Proc. IEEE Intl. Conf. on Data Mining, 2013).(w/ Z. Jiang et al.)

# Spatial Decision Tree

## Traditional decision tree

Inputs: table of records

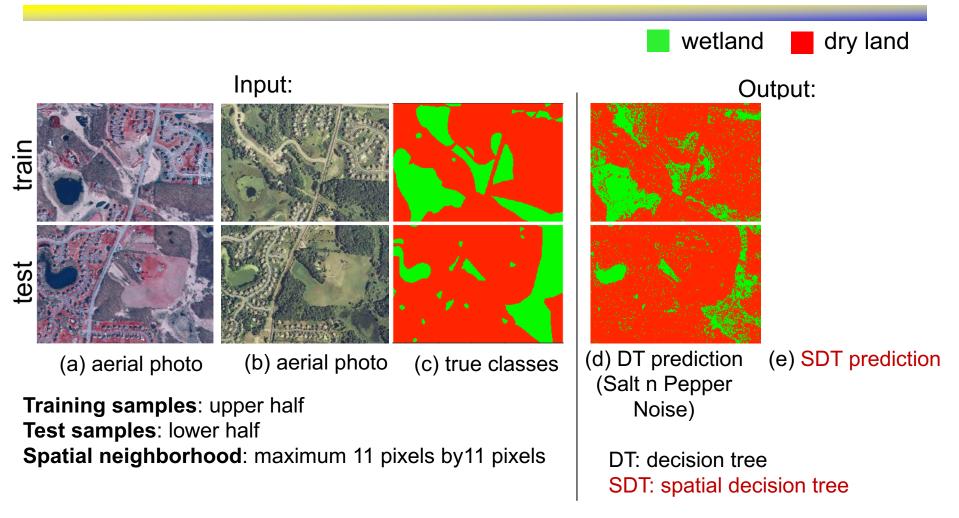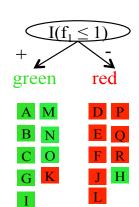| ID | $f_1$ | $f_2$ | $\Gamma_1$ | class |
|----|----|----|----|-------|
| A | 1 | 1 | 1 | green |
| B | 1 | 1 | 0.3 | green |
| C | 1 | 3 | 0.3 | green |
| G | 1 | 1 | 0.3 | green |
| I | 1 | 3 | 0 | green |
| K | 1 | 2 | -1 | red |
| M | 1 | 1 | 1 | green |
| N | 1 | 1 | 0.3 | green |
| O | 1 | 3 | 0.3 | green |
| D | 3 | 2 | 0.3 | red |
| E | 3 | 2 | 0.3 | red |
| F | 3 | 2 | 1 | red |
| H | 3 | 1 | -1 | green |
| J | 3 | 2 | 0 | red |
| L | 3 | 2 | 0.3 | red |
| P | 3 | 2 | 0.3 | red |
| Q | 3 | 2 | 0.3 | red |
| R | 3 | 2 | 1 | red |

$I(f_1 \leq 1)$

+  green   -  red

green: A M B N C O G K I
red: D P E Q F R J H L

Predicted map

| A | B | C | D | E | F |
| G | H | I | J | K | L |
| M | N | O | P | Q | R |

| feature test | information gain |
|--------------|------------------|
| $f_1 \leq 1$ | 0.50 |
| $f_2 \leq 1$ | 0.46 |
| $f_2 \leq 2$ | 0.19 |

## Spatial decision tree

Inputs:
- feature maps, class map
- Rook neighborhood

Feature $f_1$

| 1 | 1 | 1 | 3 | 3 | 3 |
| 1 | 3 | 1 | 3 | 1 | 3 |
| 1 | 1 | 1 | 3 | 3 | 3 |

Feature $f_2$

| 1 | 1 | 3 | 2 | 2 | 2 |
| 1 | 1 | 3 | 2 | 2 | 2 |
| 1 | 1 | 3 | 2 | 2 | 2 |

Class map

Focal function $\Gamma_1$

| 1 | .3 | .3 | .3 | .3 | 1 |
| .3 | -1 | 0 | 0 | -1 | .3 |
| 1 | .3 | .3 | .3 | .3 | 1 |

$I(f_1 \leq 1) * \Gamma_1$

+  green   -  red

green: A M B N C O G H I
red: D P E Q F R J K L

Predicted map

| A | B | C | D | E | F |
| G | H | I | J | K | L |
| M | N | O | P | Q | R |

# Location Prediction Models

- **Traditional Models,** e.g., Regression (with Logit or Probit),
  - Linear Regression, Bayes Classifier, …
- **Semi-Spatial** : auto-correlation regularizer
- Spatial Models $$\varepsilon = \|y - \beta X\|^2 + \|\beta X - \beta X_{neighbor}\|^2$$
  - Spatial autoregressive model (SAR)
  - Markov random field (MRF) based Bayesian Classifier

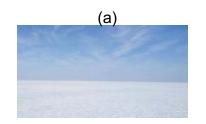| Traditional | Spatial |
|---|---|
| $y = X\beta + \varepsilon$ | $y = \rho W y + X\beta + \varepsilon$ |
| $\Pr(C_i \mid X) = \dfrac{\Pr(X \mid C_i)\,\Pr(C_i)}{\Pr(X)}$ | $\Pr(c_i \mid X, C_N) = \dfrac{\Pr(C_i)\,\Pr(X, C_N \mid c_i)}{\Pr(X, C_N)}$ |
| Neural Networks | Convolutional Neural Networks |
| Decision Trees | Spatial Decision Trees |

# Spatial Variability Challenge: Amorphous Features

**Q1.** Which images show snow ?

(a)

Runn of Kutch, Gujarat, India

(b)

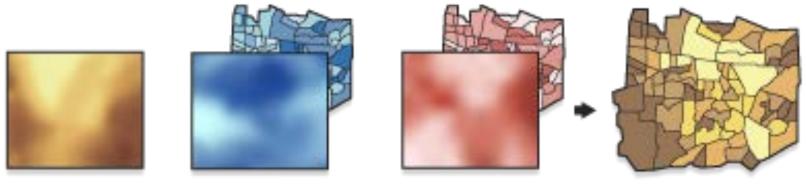White Sands, NM, USA

(c)

shutterstock.com · 116906137

Snow

**Q2.** Which geo-challenges are addressed by Convolutional Neural Network (CNN) ?
(a) High Cost of spurious and missed patterns (b) Spatial Auto-correlation
(c) Spatial Heterogeneity                                    (d) Teleconnections

# Modeling Spatial Heterogeneity: GWR

- Geographically Weighted Regression (GWR)
  - Goal: Model spatially varying relationships
  - Example: $y = X\beta' + \varepsilon'$
    Where $\beta'$ and $\varepsilon'$ are location dependent



$$\beta_0 \quad + \quad \beta_1 \text{ Population} \quad + \quad \beta_2 \text{ Income} \quad = \quad \text{Crime}$$
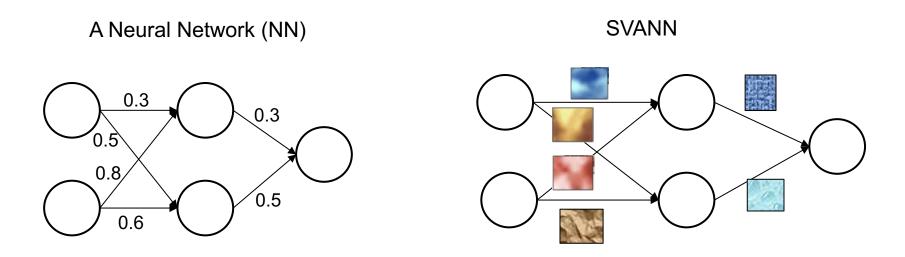
Source: resources.arcgis.com

# Spatial Variability Aware Neural Networks (SVANN)

- Each NN parameter is a map i.e., a function of location
  - Similar to Geographically Weighted Regression

A Neural Network (NN)

SVANN



- Evaluation Task:
  - Urban Garden Detection across Hennepin County, MN and Fulton County, GA.
  - SVANN outperformed OSFA by 14.34% on F1-scores.

# Spatial Challenges: Traditional Data Science

- Traditional methods not robust in face of
  - Challenge 1: Spatial continuity
  - Challenge 2: Auto-correlation, Heterogeneity , Edge-effect, …
  - Challenge 3: Noise
    - High cost of spurious patterns

**Details:** Data Science for Earth: The Earth Day Report , E. Eftelioglu, S. Shekhar, J. Hudson, L. Joppa, C. Baru, V. Janeja, et al. ACM SIGKDD Explorations Newsletter, 22(1), May 2020

# Dealing with Noise & Spurious Chance Patterns

- ## Statistics: Deal with Noise
  - Quantify uncertainty, confidence, …
  - Is it (statistically) significant?
  - Is it different from a chance event or rest of dataset?
    - e.g., SaTScan finds circular hot-spots

- ## Spatial Statistics, Spatial Data Mining
  - Auto-correlation, Heterogeneity, Edge-effect, …



■ Pump sites
☆ Deaths from cholera

Number of cases: 144
Expected cases: 62.13
Log likelihood ratio: 60.37
P-value: 0.001

Soho

Input: 250 cholera cases (multiple fatalities are simplified as a single case.)



**SaTScan**™
Software for the spatial, temporal, and space-time scan statistics

# Spatial Scan Statistics (SatScan)

- Goal: Omit chance clusters

- Ideas: Likelihood Ratio, Statistical Significance

- Steps
    - Enumerate candidate zones & choose zone X with highest likelihood ratio (LR)
        - $LR(X) = p(H1|data) / p(H0|data)$
        - H0: points in zone X show complete spatial randomness (CSR)
        - H1: points in zone X are clustered

    - If $LR(Z) >> 1$ then test statistical significance
        - Check how often is $LR( CSR ) > LR(Z)$
            using 1000 Monte Carlo simulations

# Beyond SatScan: Spatial Concept/Theory-Aware Hotspots

- Geographic features, e.g., rivers, streams, roads, …
  - Hot-spots => Hot Geographic-features, e.g., Linear Hotspots
- Spatial Theories, e.g,, environmental criminology
  - Circles ➜ Doughnut holes



Pedestrian fatalities
Orlando, FL

Circular hotspots
by SatScan

Linear hotspots

**Details:** Significant Linear Hotspot Discovery, IEEE Transactions on Big Data, 3(2):140-153, 2017.
(Summary in Proc. Geographic Info. Sc., Springer LNCS 8728, pp. 284-300, 2014.)

# Hotel That Enlivened the Bronx Is Now a 'Hot Spot' for Legionnaires'

By WINNIE HU and NOAH REMNICK    AUG. 10, 2015

## Contaminated Cooling Towers

Five buildings have been identified as the potential source of the Legionnaires' disease outbreak in the South Bronx.

■ Possible sources of Legionnaires' outbreak
■ Additional sites found with legionella bacteria
● Locations of people with Legionnaires'

Verizon office building

Harlem River

E 167TH

BRONX

E 161ST ST

Concourse Plaza

Lincoln Medical Center

Opera House Hotel

Streamline Plastics Company

East River

2,000 Feet

Source: New York Mayor's Office

By The New York Times

The Opera House Hotel is at the center of the outbreak. Edwin J. Torres for The New York Times

# Legionnaires' Disease Outbreak in New York



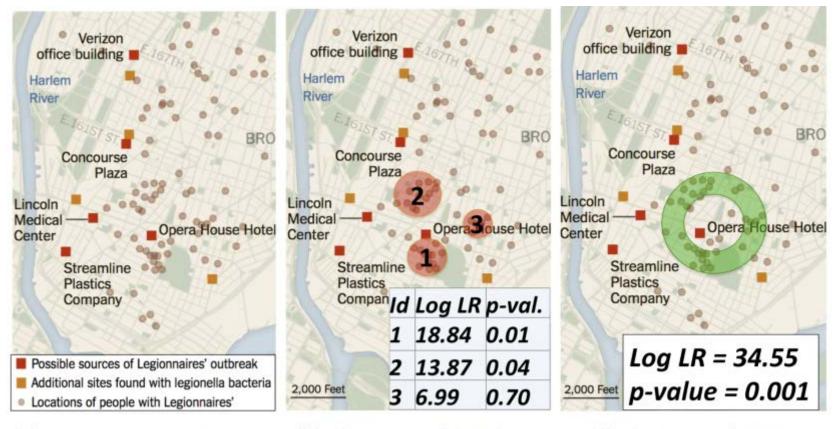(a) Legionnaire's in New York (2015)

(b) Output of SaTScan

| Id | Log LR | p-val. |
|----|--------|--------|
| 1 | 18.84 | 0.01 |
| 2 | 13.87 | 0.04 |
| 3 | 6.99 | 0.70 |

(c) Output of RHD

*Log LR = 34.55*
*p-value = 0.001*

Legend:
- ■ Possible sources of Legionnaires' outbreak
- ■ Additional sites found with legionella bacteria
- • Locations of people with Legionnaires'

Details: Ring-Shaped Hotspot Detection, IEEE Trans. Know. & Data Eng., 28(12), 2016.
(A Summary in Proc. IEEE ICDM 2014) (w/ E. Eftelioglu et al.)
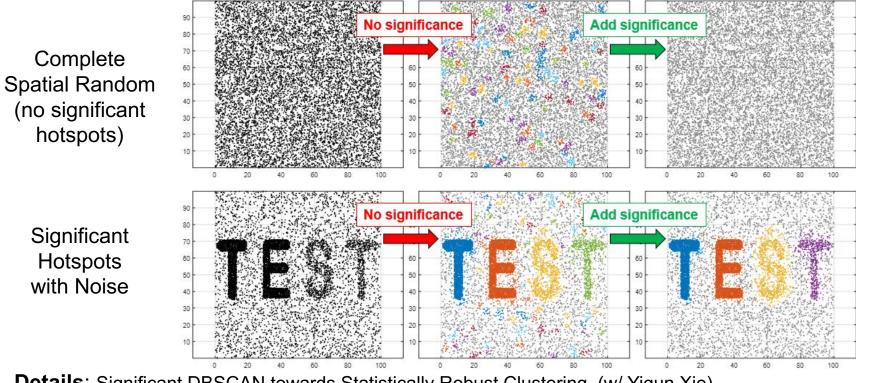
# Robust Clustering (Hotspot Detection)

- **Problem definition**
  - **Inputs:** Collection of event locations, Test statistic; Significance level
  - **Output:** Significant clusters (hotspots)
  - **Constraints:** Avoid chance patterns despite non-trivial noise in data

- **Limitations of Related Work**
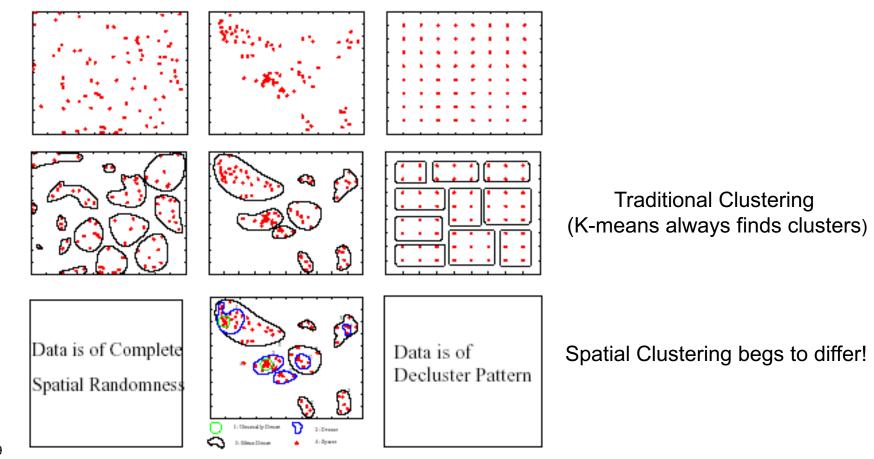  - DBSCAN cannot avoid chance patterns
  - SaTScan cannot detect clusters of arbitrary shapes
- **Contributions**
  - Significance modeling in DBSCAN
  - A fast dual-convergence algorithm



**Details**: Significant DBSCAN towards Statistically Robust Clustering, (w/ Yiqun Xie),
In Proc. 16th Intl. Symposium on Spatial and Temporal Databases (SSTD), 2019, ACM. **(Best Paper Award)**

# Limitation of Traditional Clustering

- Challenge: One size does not fit all
  - Prediction error vs. model bias, Cost of false positives, …
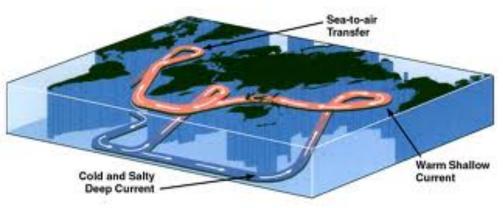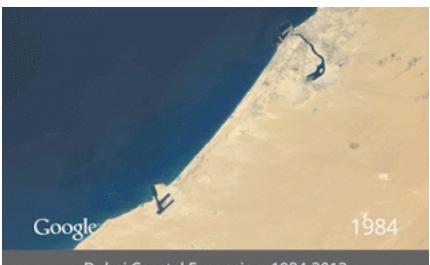- Example. Clustering: Find groups of points



Traditional Clustering
(K-means always finds clusters)

Spatial Clustering begs to differ!

# What has changed? **Spatial Data Revolution**

| Spatial | Last Century | Last Decade |
|---|---|---|
| **Data** | Few satellites and sensors | Nano-satellites, Billions of GPS enabled smartphones |
| **Data Access** | Need special hardware and network | Cloud based repositories, e.g., Earth on AWS |
| **Spatial Platforms** | ESRI Arc/Info | SQL3/OGC, e.g., Postgis, ESRI GIS Tools for Hadoop, Google Earth Engine |
| **Spatial Data Science** | Spatial Patterns, e.g., hotspots (SatScan, ESRI Geostatistics) | (a) Spatial Network Patterns, e.g., linear hotspots<br>(b) Spatio-temporal (ST) patterns, e.g., Change time-series (Google Timelapse) |
| **Spatial Visualization** | Quilt: MS Terraserver<br>Fly through: Google Earth | (a) Space time: Timelapse<br>(b) There Dimensions |

# Towards Time-Travel and Depth in Virtual Globes

- Virtual globes are snapshots

- How to add time? depth?
  - Ex. Google Earth Engine, NASA NEX
  - Ex. Google Timelapse: 260,000 CPU core-hours for global 29-frame video
  - https://earthengine.google.com/timelapse/
  - Salt Lake, Bidhannagar, Kolkata, WB, India
  - UMN, Minneapolis;airport, MN,USA

- How may one convey provenance,



Dubai Coastal Expansion, 1984-2012





Sea-to-air Transfer

Cold and Salty Deep Current

Warm Shallow Current

# A UCGIS Call to Action:
# Bringing the Geospatial Perspective to Data Science Degrees and Curricula

Data that are geographically referenced or contain some type of location markers are both common and of high value (e.g., data subject to state-specific policies, laws and regulations; demographic data from the census; location traces of smartphones and vehicles; remotely sensed imagery from satellites, aircraft and small unmanned aerial vehicles; volunteered geographic information; geographically referenced social media postings). A 2011 McKinsey Global Institute report estimates a value of "about $600 billion annually by 2020" from leveraging personal location data[2] to reduce fuel waste, improve health outcomes, and better match products to consumer needs. Spatial data are critical for societal priorities such as national security, public health & safety, food, energy, water, smart cities, transportation, climate, weather, and the environment. For example, remotely-sensed satellite imagery is used to monitor not only weather and climate but also global crops[3] for early warnings and planning to avoid food shortages.

University Consortium for
**GEOGRAPHIC INFORMATION SCIENCE**

*Summer 2018*

# One Size Data Science Does not Fit All Data!

However, spatial data presents unique data science challenges. Recent court cases that address gerrymandering, the manipulation of geographic boundaries to favor a political party, offer a high-profile example. Instances of such exploitation of the modifiable areal unit problem (or dilemma) is not limited to elections since the MAUP affects almost all traditional data science methods in which results (e.g., correlations) change dramatically by varying geographic boundaries of spatial partitions. The fundamental geographic qualities of spatial autocorrelation, which assumes properties of geographically proximate places to be similar, and geographic heterogeneity, where no two places on Earth are exactly alike, violate assumptions of sample independence and randomness that underlie many conventional statistical methods. Other spatial challenges include how to choose between a plurality of projections and coordinate systems and how to deal with the imprecision, inaccuracy, and uncertainty of location

## A UCGIS Call to Action:
## Bringing the Geospatial Perspective to Data Science Degrees and Curricula

University Consortium for
**GEOGRAPHIC INFORMATION SCIENCE**

*Summer 2018*

# Spatial Data Science Tools

measurements. To deal with such challenges, practitioners in many fields including agriculture, weather forecast, mining, and environmental science incorporate *geospatial data science*[4] methods such as spatially-explicit models, spatial statistics[5], geo-statistics, geographic data mining[6], spatial databases[7], etc.

[4] Y. Xie et al., Transdisciplinary Foundations of Geospatial Data Science, *ISPRS Intl. Jr. of Geo-Informatics*, 6(12):395-418, 2017. DOI: 10.3390/ijgi6120395.

[5] N. Cressie, *Statistics for Spatial Data*, Wiley, 1993 (1st ed.), 2015 (Revised ed.).

[6] H. Miller and J. Han, *Geographic Data Mining and Knowledge Discovery*, CRC Press, 2009 (2nd Ed.).

[7] S. Shekhar and S. Chawla, *Spatial Databases: A Tour*, Prentice Hall, 2003.

## A UCGIS Call to Action:
## Bringing the Geospatial Perspective to Data Science Degrees and Curricula

University Consortium for
**GEOGRAPHIC INFORMATION SCIENCE**

*Summer 2018*

# **Summary :** One size data science does not fit all

- Spatial Data are ubiquitous & important

- Traditional Data Science Tools are inadequate
    – Gerrymandering, Spatial Auto-correlation, …

- **Ask:**
    – Spatial Data Science Methods
    – Spatial Statistics, Spatial Data Mining, SDBMS, …

One size does NOT fit all.

# References : Surveys, Overviews

- Spatial Computing, MIT Press (Essential Knowledge Series), 2020. (ISBN: 9780262538046).
- Spatial Computing ( html , short video , tweet ), Communications of the ACM, 59(1):72-81, January, 2016.
- Transdisciplinary Foundations of Geospatial Data Science ( html , pdf ), ISPRS Intl. Jr. of Geo-Informatics, 6(12):395-429, 2017. ( doi:10.3390/ijgi6120395 )
- Spatiotemporal Data Mining: A Computational Perspective , ISPRS Intl. Jr. on Geo-Information, 4(4):2306-2338, 2015 (DOI: 10.3390/ijgi4042306).
- Identifying patterns in spatial information: a survey of methods ( pdf ), Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 1(3):193-214, May/June 2011. (DOI: 10.1002/widm.25).
- Theory-Guided Data Science: A New Paradigm for Scientific Discovery from Data, IEEE Transactions on Knowledge and Dat Mining, 29(10):2318-2331, June 2017. ( DOI: 10.1109/TKDE.2017.2720168 ).
- Parallel Processing over Spatial-Temporal Datasets from Geo, Bio, Climate and Social Science Communities: A Research Roadmap. IEEE BigData Congress 2017: 232-250.
- Spatial Databases: Accomplishments and Research Needs, IEEE Transactions on Knowledge and Data Engineering, 11(1):45-55, 1999.