# University of Minnesota

# Machine Learning for Big Spatial Data and Applications

Mohamed F. Mokbel
University of Minnesota

# The Ubiquity of Big Spatial Data and Applications



4 TB every day

50PB / year

UBER

Share it

NASA

OpenStreetMap

yelp

foursquare

flickr

# Big Spatial Data in Agriculture

## Spatial Data and Precision Agriculture

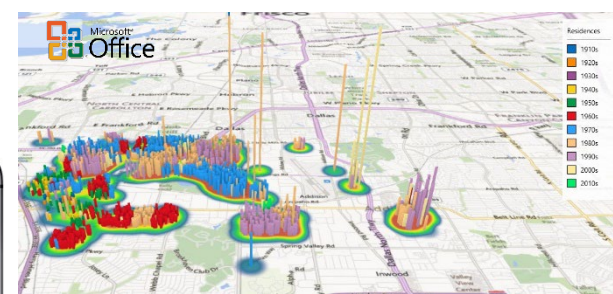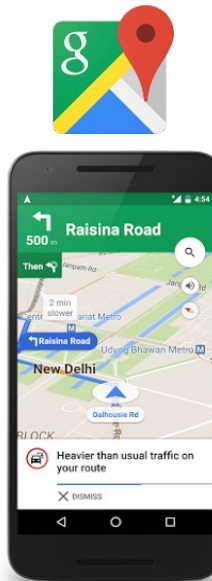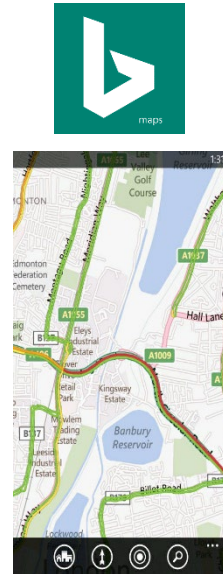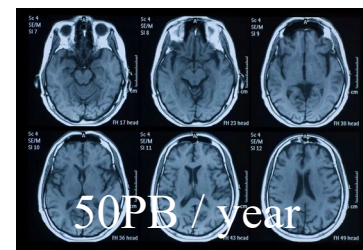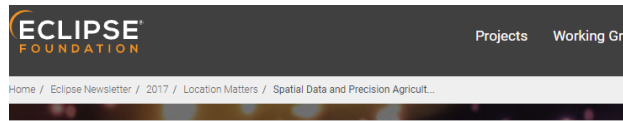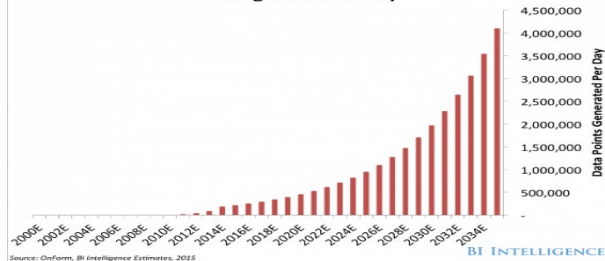Precision Agriculture is a methodology of farm management that relies on data, and data analysis to support the farmer's decision-making process to decrease inputs.

The origin of Precision Agriculture begins with researchers collecting soil samples, and using spatial statistics methods to determine the different soil types in a field. From this analysis, the researchers developed soil maps. Farms were early adopters of both GPS and Geographic Information Systems (GIS) technologies. As civilian GPS became more accurate, farms started to utilize this technology to increase the accuracy of operational spatial data. Collecting spatial data from equipment and sensors that allowed farms to pinpoint the high yield areas. Also, using GPS data to determine where to increase or decrease pesticides, fertilizers use and irrigation.

**Average Farm Per Day**

Source: OnFarm, BI Intelligence Estimates, 2015

**BI INTELLIGENCE**

### Agriculture Technology: How GIS Can Help You Win the Farm

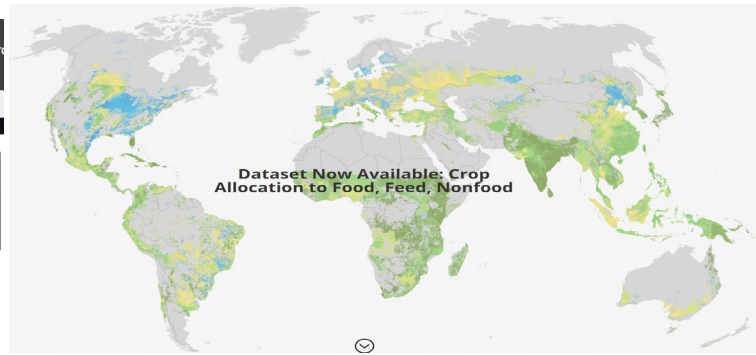By: GISGeography • Last Updated: August 4, 2021

### Agriculture Technology from Location

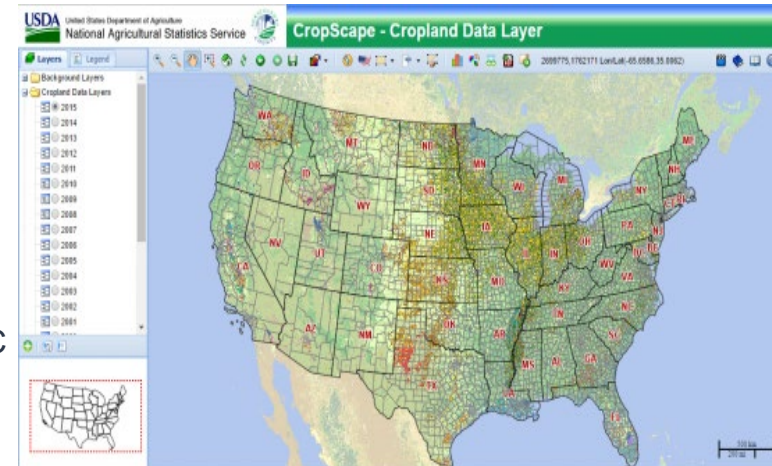Today's farmers use sophisticated **agriculture technology** because they can save time and money.

Because crops are location-based, this makes Geographic Information Systems (GIS) an EXTREMELY relevant tool for farmers.

For example, farmers use precision GPS on the field to save fertilizer. Also, satellites and drones collect vegetation, topography, and weather information from the sky.

**Dataset Now Available: Crop Allocation to Food, Feed, Nonfood**

**EARTHSTAT**

EarthStat serves geographic data sets that help solve the grand challenge of feeding a growing global population while reducing agriculture's impact on the environment.

CropScape is developed by USDA-NASS where farmers can see *what* crops are growing *where* and *how much*. CropScape is also used for food security, land-cover change and pesticide control: https://nassgeodata.gmu.edu/CropScape/
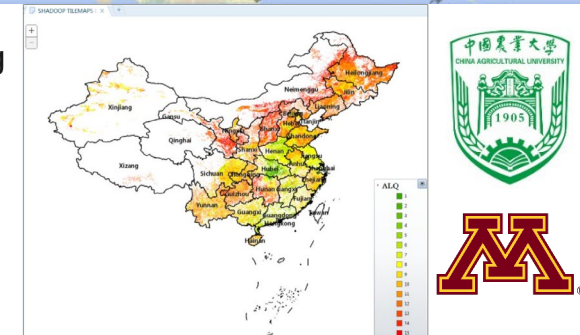
**CropScape - Cropland Data Layer**

### LandQ$^{v2}$: A MapReduce-Based System for Processing Arable Land Quality Big Data

by Xiaochuang Yao [1,*], Mohamed F. Mokbel [2], Sijing Ye [3], Guoqing Li [1], Louai Alarabi [2], Ahmed Eldawy [4], Zuliang Zhao [5], Long Zhao [5] and Dehai Zhu [5,*]

[1] Institute of Remote Sensing and Digital Earth, Chinese Academy of Sciences, Beijing 100094, China

[2] Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN 55455, USA

[3] State Key Laboratory of Earth Surface Processes and Resource Ecology, Beijing Normal University, Beijing 100875, China

[4] Department of Computer Science and Engineering, University of California, Riverside, CA 92521, USA

[5] College of Information and Electrical Engineering, China Agricultural University, Beijing 100083, China

# Big Spatial Data in Transportation



Routing & Scheduling Optimisation

# Big Spatial Data in Polar Regions



National Science Foundation
WHERE DISCOVERIES BEGIN

SEARCH

RESEARCH AREAS    FUNDING    AWARDS    DOCUMENT LIBRARY    NEWS    ABOUT NSF

Awards

Search Awards

Recent Awards

Presidential and Honorary
Awards

Award Abstract # 2118285
HDR Institute: HARP- Harnessing Data and Model Revolution in the Polar Regions

| NSF Org: | OAC<br>Office of Advanced Cyberinfrastructure (OAC) |
| Awardee: | UNIVERSITY OF MARYLAND BALTIMORE COUNTY |
| Initial Amendment Date: | September 15, 2021 |

UMBC · TEXAS The University of Texas at Austin · UAF UNIVERSITY OF ALASKA FAIRBANKS · Boulder · AMHERST COLLEGE · BOWIE STATE UNIVERSITY 1865 · NASA · DARTMOUTH

The HDR Institute aims to harness massive heterogeneous, noisy, and discontinuous data in space and time and integrate data with numerical and physical models

Researchers at i-HARP are investigating novel data science techniques including deep generative adversarial networks, graph neural networks, meta learning, hybrid networks, physics-informed machine learning, causal artificial intelligence, data assimilation, spatio-temporal deep learning, and scalable algorithms.

# Big Spatial Data for .....

**Remote Sensing**

Jan-2009

https://lpdaac.usgs.gov
LP DAAC archive exceeds 1PB
5 Trillion points Temperature data
Vegetation data at 250m2
resolution (16 times larger)

GIS INNOVATION CENTER

72 months × 14 Billion points/month = **1 Trillion points**

A. Eldawy, M. Mokbel, S. Alharthi, A. Alzaidy, K. Tarek, S. Ghani. "*SHAHED: A MapReduce-based System for Querying and Visualizing Spatio-temporal Satellite Data*". **ICDE 2015**

**Telco Data**

**Telco Big Data Awareness**

The aim of this project is to develop next generation telco big data management architectures that can help in understanding urban phenomena (e.g., traffic in a city, mobility patterns for emergency response or city planning, improve the Quality of Service) at a very high spatio-temporal resolution. The project deals with algorithms and structures to ingest in the most compact manner huge amounts of network logs perform big data exploration and analytics within a tolerable elapsed time.

CITY — collect

TELCO
Big Data
Analytics
Knowledge
Open Data

optimize

SMART CITY ENABLERS
Municipality/ City Planning
Public Services/ Utilities
Startup/ Smart App
Transportation/ Traffic
Company/ Marketing

C. Costa, G. Chatzimilioudis, D. Zeinalipour-Yazti, M. Mokbel: Efficient Exploration of Telco Big Data with Compression and Decaying. ICDE 2017: 1332-1343

University of Cyprus

. . . . . . .

# Meanwhile, the Rise of Machine Learning



*"Machine learning is a core, transformative way by which we're rethinking everything we're doing."*
*-Google CEO Sundar Pichai*

# Machine Learning meets Big Spatial Data

# Knowledge Base Construction



Probabilistic Knowledge Base Construction System

**Relations Extraction**

**Factual Scores Inference**

Knowledge Base Rules

tables, relations of facts

Spouses KB

| Person 1 | Person 2 |
|----------|----------|
| Barack | Michelle |

NELL

yago select knowledge

SystemT

StatSnowBall

DARPA MEMEX

Fight Human Trafficking Crime Investigation

DeepDive

liXto DELIVERING COMPETITIVE ADVANTAGE

DBLife

Google Vault

ProbKB

SCIENTIFIC AMERICAN™ WSJ 60 MINUTES CNN Forbes BBC

appleinsider

Apple acquires "dark data" specialist Lattice Data for $200M

By Daniel Eran Dilger
Saturday, May 13, 2017, 12:29 pm PT (03:29 pm ET)

# DeepDive with Spatial Data …

**Crime rates in Minnesota**

Crimes

| City | C | E |
|------|---|-----|
| Minneapolis | 1 | 0.7 |
| St. Paul | ? | 0.7 |
| Eagan | ? | 0.7 |
| Rochester | ? | 0.7 |

Education

**Data**

Minneapolis
St. Paul
Eagan
Rochester

```
P1: City X has high crime rate
P2: Cities X&Y have same education level



Rule: P1&P2 ➔ Y has high crime rate
```

**Inference Rules**

DeepDive

| City | Confidence | | |
|------|-----------|---|---|
| St. Paul | 0.5 | | |
| Eagan | 0.5 | | |
| Rochester | 0.5 | | |

**Result**

# DeepDive with Spatial Data …

**Crime rates in Minnesota**

Crimes

| City | C | E |
|------|---|-----|
| Minneapolis | 1 | 0.7 |
| St. Paul | ? | 0.7 |
| Eagan | ? | 0.7 |
| Rochester | ? | 0.7 |

Education

**Data**

Minneapolis
St. Paul
Eagan
Rochester

```
P1: City X has high crime rate
P2: Cities X&Y have same education level
P3: Cities X&Y are within 80 miles


Rule: P1&P2  ➔ Y has high crime rate
Rule: P1&P2&P3 ➔ Y has high crime rate
```

**Inference Rules**

**DeepDive**

**Result**

| City | Confidence | |
|------|-----|-----|
| St. Paul | ~~0.5~~ | 0.7 |
| Eagan | ~~0.5~~ | 0.7 |
| Rochester | ~~0.5~~ | 0 |

# DeepDive with Spatial Data …

## Crime rates in Minnesota

Crimes

| City | C | E |
|------|---|-----|
| Minneapolis | 1 | 0.7 |
| St. Paul | ? | 0.7 |
| Eagan | ? | 0.7 |
| Rochester | ? | 0.7 |

Education

**Data**

Minneapolis
St. Paul
Eagan
Rochester

```
P1: City X has high crime rate
P2: Cities X&Y have same education level
P3: Cities X&Y are within 80 miles


Rule: P1&P2 ➔ Y has high crime rate
Rule: P1&P2&P3 ➔ Y has high crime rate
```

**Inference Rules**

## Ebola infection rates in Liberia

Infections

| County | I | S |
|--------|---|-----|
| Montserrado | 1 | 0.6 |
| Margibi | ? | 0.6 |
| Bong | ? | 0.6 |
| Gbarpolu | ? | 0.6 |

Sanitation

Gbarpolu
Bong
Margibi
Montserrado
150 miles

```
P1: County X has high Ebola infection rate
P2: Counties X&Y have same sanitation level



Rule: P1&P2 ➔ Y has high infection rate
```

**DeepDive**

**Result**

**DeepDive**

| City | Confidence | |
|------|-----------|---|
| St. Paul | 0.5 | 0.7 |
| Eagan | 0.5 | 0.7 |
| Rochester | 0.5 | 0 |

| City | Confidence | |
|------|-----------|---|
| Margibi | 0.54 | |
| Bong | 0.52 | |
| Gbarpolu | 0.63 | |

# DeepDive with Spatial Data …

## Crime rates in Minnesota

Crimes

| City | C | E |
|------|---|-----|
| Minneapolis | 1 | 0.7 |
| St. Paul | ? | 0.7 |
| Eagan | ? | 0.7 |
| Rochester | ? | 0.7 |

Education

Minneapolis
St. Paul
Eagan
Rochester

**Data**

```
P1: City X has high crime rate
P2: Cities X&Y have same education level
P3: Cities X&Y are within 80 miles


Rule: P1&P2 ➔ Y has high crime rate
Rule: P1&P2&P3 ➔ Y has high crime rate
```
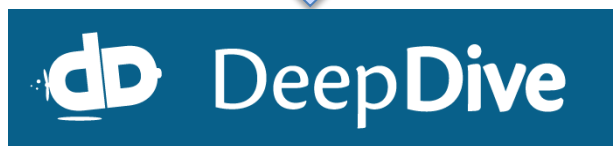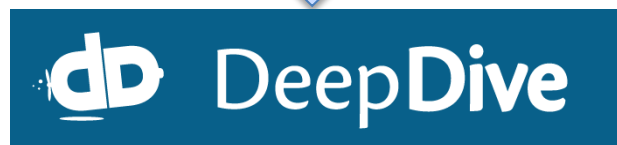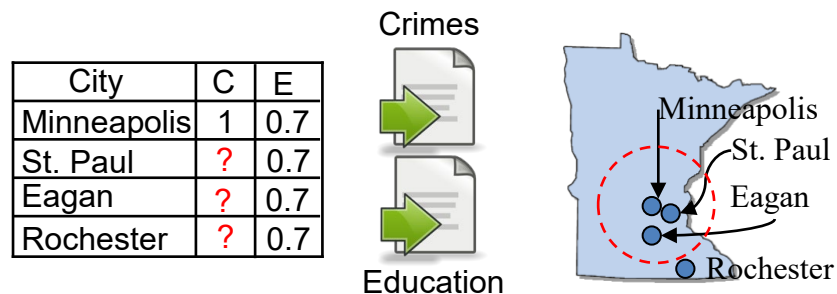
**Inference Rules**

## Ebola infection rates in Liberia

Infections

| County | I | S |
|--------|---|-----|
| Montserrado | 1 | 0.6 |
| Margibi | ? | 0.6 |
| Bong | ? | 0.6 |
| Gbarpolu | ? | 0.6 |

Sanitation

Gbarpolu
Bong
Margibi
Montserrado
150 miles

```
P1: County X has high Ebola infection rate
P2: Counties X&Y have same sanitation level
P3: Counties X&Y are within 150 miles


Rule: P1&P2 ➔ Y has high infection rate
Rule: P1&P2&P3 ➔ Y has high infection rate
```

**DeepDive**

**DeepDive**

**Result**

| City | Confidence | |
|------|------|------|
| St. Paul | ~~0.5~~ | 0.7 |
| Eagan | ~~0.5~~ | 0.7 |
| Rochester | ~~0.5~~ | 0 |

| City | Confidence | |
|------|------|------|
| Margibi | ~~0.54~~ | 0.51 |
| Bong | ~~0.52~~ | 0.45 |
| Gbarpolu | ~~0.63~~ | 0.06 |

# DeepDive with Spatial Data …

## Crime rates in Minnesota

Crimes

| City | C | E |
|------|---|---|
| Minneapolis | 1 | 0.7 |
| St. Paul | ? | 0.7 |
| Eagan | ? | 0.7 |
| Rochester | ? | 0.7 |

Education

Minneapolis
St. Paul
Eagan
Rochester

**Data**

## Ebola infection rates in Liberia

Infections

| County | I | S |
|--------|---|---|
| Montserrado | 1 | 0.6 |
| Margibi | ? | 0.6 |
| Bong | ? | 0.6 |
| Gbarpolu | ? | 0.6 |

Sanitation

Gbarpolu
Bong
Margibi
Montserrado

```
P1: City X has high crime rate
P2: Cities X&Y have same education level
P3: Cities X&Y are within 80 miles
P3: The closer Y&X the higher Y crime rate

Rule: P1&P2  ➔  Y has high crime rate
Rule: P1&P2&P3  ➔  Y has high crime rate
```
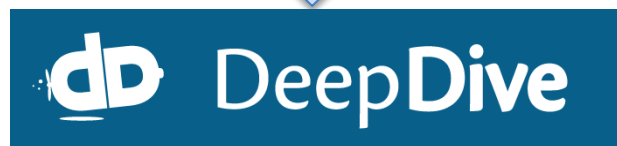
**Inference Rules**

```
P1: County X has high Ebola infection rate
P2: Counties X&Y have same sanitation level
P3: Counties X&Y are within 150 miles
P3: The closer Y&X the higher Y infect rate

Rule: P1&P2  ➔  Y has high infection rate
Rule: P1&P2&P3  ➔  Y has high infection rate
```
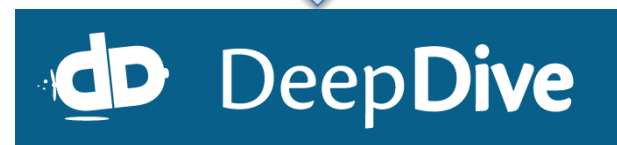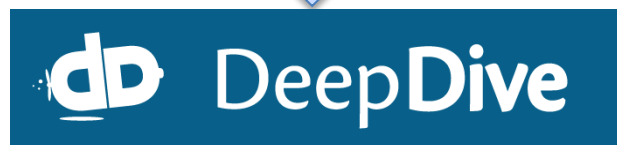
**DeepDive**

**DeepDive**

**Result**

| City | Confidence | |
|------|-----------|---|
| St. Paul | 0.5 | 0.7 |
| Eagan | 0.5 | 0.7 |
| Rochester | 0.5 | 0 |

| City | Confidence | |
|------|-----------|---|
| Margibi | 0.54 | 0.51 |
| Bong | 0.52 | 0.45 |
| Gbarpolu | 0.63 | 0.06 |

# Where Is the Problem?

- DeepDive is built on top of Markov Logic Networks (MLN)
  - MLN is designed for *binary logic* only
    - E.g., bitwise-AND, bitwise-OR, and imply

- MLN is not spatially- aware
  - It can not interpret the *gradual semantics* of spatial predicates
    - E.g., P3: The closer Y&X the higher Y infect rate

Need to build *Spatial Markov Logic Networks (SMLN)*, a full-fledged MLN framework with a native support for spatial data and applications

# Markov Logic Networks (MLN)

Machine learning — a form of artificial intelligence that uses algorithms and large data sets to derive insights in real time — is way more than hype.

Gartner predicts that by 2018, 45 percent of the fastest–growing companies will have fewer employees than instances of smart machines.

It's clear that machine learning offers companies a competitive advantage, but is it something that small– and medium–sized business can adopt? The algorithms churning the data are often opaque, and things can go wrong, from the humorous (automated email replies that write "I love you" to a

**Need experts and highly-trained scientists, specially for deep learning**

☐ **MLN is an end-to-end ML solution**
- ☐ Covers wide range of ML problems
- ☐ Thousands of lines of ML code can be done in few MLN formulas

Rules as MLN formulas → **Markov Logic Network (MLN)** → Rule weights →

**Alchemy - Open Source AI**

ACM SIGMOD/PODS International Conference on Management of Data
June 10 – June 15, 2018    Houston, TX, USA
SIGMOD 2018: Keynote Talks
Machine Learning for Data Management: Problems and Solutions

*Tuffy*

Scalable RDBMS-based MLN System

**DeepDive**

DARPA MEMEX

HoloClean

# MLN Architecture

F₁: Illiteracy → Crime [0.5]
F₂: Crime ^ Non-safety [0.7]

**Application Developer**

**System Admin**

**Applications (e.g., DeepDive, … )**

Configs

**Language**

Application Rules

**Propositional Logic Language**

**Inference**

**Gibbs Sampling Algorithm**

**Learning**

**Gradient Descent Optimization**

Inference Output

Inference Iterations

Read / Update

Compiled Rules

Factor Graph Correlations

Learned Correlations' Weights

**Grounding**

**In-DBMS Factor Graph Construction**

Input and Supervision Data

Inferred Variables' Values

**In-memory Factor Graph Index**

Factor Graph

0.5 F₁ F₂ 0.7

I C N

Factor Graph Variables

# Spatial MLN Architecture

$F_1$: Illiteracy → Crime [0.5]
$F_2$: Crime ∧ Non-safety [0.7]

**Application Developer**

**System Admin**

**Configs**

**Language**

**DDlog Language with Spatial Extensions**

**Application Rules**

**Applications (e.g., Sya, Flash, … )**

**Inference**

**In-memory Spatial Factor Graph Index**

**Learning**

**Spatial Gradient Descent Optimization**

**Inference Output**

**Compiled Rules**

**Inference Iterations**

**Read / Update**

**Factor Graph Correlations**

**Learned Correlations' Weights**

**Grounding**

**In-DBMS Spatial Factor Graph Construction**

**Input and Supervision Data**

**Inferred Variables' Values**

**In-memory Spatial Factor Graph Index**

**Factor Graph**  0.5 $F_1$  $F_2$ 0.7
I  C  N

**Factor Graph Variables**

# Sya

## Crime rates in Minnesota

Crimes

| City | C | R |
|------|---|---|
| Minneapolis | 1 | 0.7 |
| St. Paul | ? | 0.7 |
| Eagan | ? | 0.7 |
| Rochester | ? | 0.7 |

Education



Minneapolis
St. Paul
Eagan
Rochester

**Data**

## Ebola infection rates in Liberia

Infections

| County | I | S |
|--------|---|---|
| Montserrado | 1 | 0.6 |
| Margibi | ? | 0.6 |
| Bong | ? | 0.6 |
| Gbarpolu | ? | 0.6 |

Sanitation



```
P1: City X has high crime rate
P2: Cities X&Y have same education level
P3: Cities X&Y are within 80 miles
P3: The closer Y&X the higher Y crime rate

Rule: P1&P2 ➜ Y has high crime rate
Rule: P1&P2&P3 ➜ Y has high crime rate
```

**Inference Rules**

```
P1: County X has high Ebola infection rate
P2: Counties X&Y have same sanitation level
P3: Counties X&Y are within 150 miles
P3: The closer Y&X the higher Y infect rate

Rule: P1&P2 ➜ Y has high infection rate
Rule: P1&P2&P3 ➜ Y has high infection rate
```

## Sya

## Sya

**Result**

| City | Confidence | | |
|------|-----|-----|-----|
| St. Paul | 0.5 | 0.7 | 0.9 |
| Eagan | 0.5 | 0.7 | 0.7 |
| Rochester | 0.5 | 0 | 0.3 |

| City | Confidence | | |
|------|------|------|------|
| Margibi | 0.54 | 0.51 | 0.76 |
| Bong | 0.52 | 0.45 | 0.53 |
| Gbarpolu | 0.63 | 0.06 | 0.22 |

# Machine Learning meets Big Spatial Data

# Spatial MLN Architecture



F₁: Illiteracy → Crime [0.5]
F₂: Crime ^ Non-safety [0.7]

**Application Developer**

**System Admin**

**Language**
DDlog Language with **Spatial** Extensions

**Application Rules**

**Applications** (e.g., **Sya**, Flash, … )

**Configs**

**Inference**
In-memory **Spatial** Factor Graph Index

**Learning**
**Spatial** Gradient Descent Optimization

**Inference Output**

**Inference Iterations**

**Read / Update**

**Compiled Rules**

**Factor Graph Correlations**

**Learned Correlations' Weights**

**Grounding**
In-DBMS **Spatial** Factor Graph Construction

**Input and Supervision Data**

**Inferred Variables' Values**

**In-memory Spatial** Factor Graph Index

**Factor Graph**  0.5 F₁  F₂ 0.7
I  C  N

**Factor Graph Variables**

# Spatial (Autologisitc) Regression

■ Find whether a spatial phenomenon exists or not, based on neighbor values and features
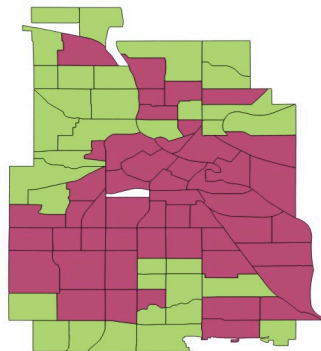
**Weather Prediction**

**Birds Migration**

$(x_1 = 0, x_2 = 1)$    $(x_1 = 1, x_2 = 1)$

Features

| $I_1$ 0 | $I_2$ 1 | $I_3$ 1 | $I_4$ 1 |
|---|---|---|---|
| $I_5$ 1 | $I_6$ 1 | $I_7$ 0 | $I_8$ 1 |
| $I_9$ 1 | $I_{10}$ 1 | $I_{11}$ 0 | $I_{12}$ 0 |
| $I_{13}$ 1 | $I_{14}$ ? | $I_{15}$ 0 | $I_{16}$ 0 |

Phenomenon value

$(x_1 = 1, x_2 = 0)$

**Land Cover**

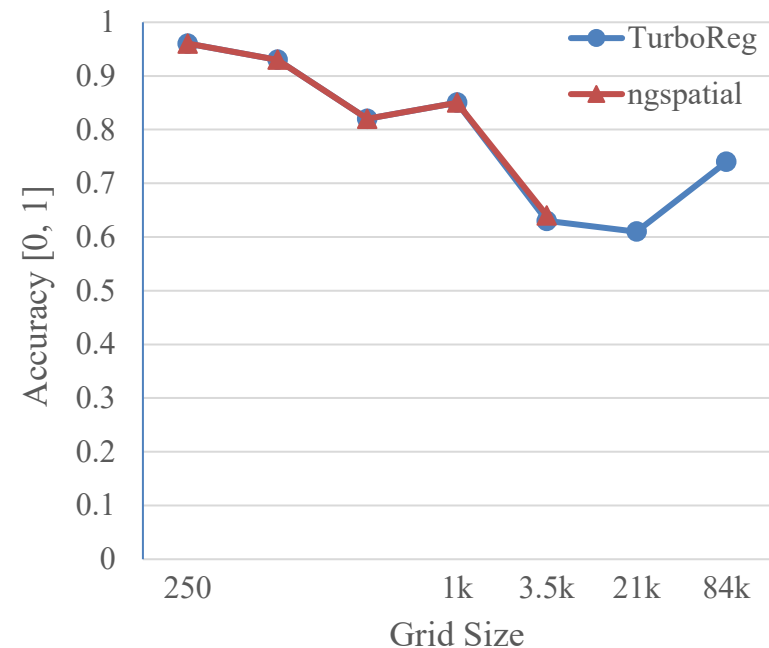**Crimes Distribution**

**Missing value**

**Features**

**Neighbor values**

$$\log \frac{Pr(z_i = 1 | \mathcal{X}, \mathcal{Z}_{\mathcal{N}_i})}{Pr(z_i = 0 | \mathcal{X}, \mathcal{Z}_{\mathcal{N}_i})} =$$

$$\sum_{j=1}^{m} \beta_j x_j + \eta \sum_{k \in \mathcal{N}_i} z_k$$

**Regression Parameters**

Learning regression parameters for 80K cells takes more than one day ☹

# Spatial Regression as SMLN Problem

$$\log \frac{Pr(z_i=1|\mathcal{X},\mathcal{Z}_{\mathcal{N}_i})}{Pr(z_i=0|\mathcal{X},\mathcal{Z}_{\mathcal{N}_i})} = \sum_{j=1}^{m} \beta_j x_j + \eta \sum_{k \in \mathcal{N}_i} z_k$$

Regression Equation → **SMLN Transformation** → SMLN Rules → **SMLN Engine** → Rule weights **=** Regression Parameters

$$\log \frac{Pr(z_i=1|\mathcal{X},\mathcal{Z}_{\mathcal{N}_i})}{Pr(z_i=0|\mathcal{X},\mathcal{Z}_{\mathcal{N}_i})} = \beta_1 x_1 + \eta \sum_{k \in \mathcal{N}_i} z_k$$

→ **SMLN Transformation** → SMLN Rules → **SMLN Engine** → **ß₁ , η**

$$\log \frac{Pr(z_1=1)}{Pr(z_1=0)} = \beta_1 x_1 + \eta z_2 + \eta z_3$$

$$\vdots$$

$$\log \frac{Pr(z_4=1)}{Pr(z_4=0)} = \beta_1 x_1 + \eta z_2 + \eta z_3 + \eta z_7 + \eta z_{10}$$

$$\vdots$$

$$\log \frac{Pr(z_{16}=1)}{Pr(z_{16}=0)} = \beta_1 x_1 + \eta z_{14} + \eta z_{15}$$

| SMLN Rules |
| --- |
| [$Z_1$ ^ $X_1$, **ß₁**] |
| [$Z_1$ ^ $Z_2$, **η**] |
| [$Z_1$ ^ $Z_3$, **η**] |
| [$Z_2$ ^ $X_1$, **ß₁**] |
| [$Z_2$ ^ $Z_4$, **η**] |
| [$Z_2$ ^ $Z_5$, **η**] |
| [$Z_3$ ^ $X_1$, **ß₁**] |
| [$Z_3$ ^ $Z_4$, **η**] |
| ……. |

**Theoretical proof of the Autologistic Regression-SMLN equivalence is in the paper**

# Multinomial Autologistic Regression

■ **Prediction and feature variables are multinomial (i.e., categorical)**

❑ Domain values are predefined values (e.g., {0, 1, 2})

❑ Represent each multinomial variable with a set of binary variables

$$\text{Pivot}$$

$$\text{Prediction} \quad z_i \equiv \begin{matrix} \boxed{z_i(0)} \\ z_i(1) \\ z_i(2) \end{matrix} \quad \{0,1\}$$

$$\{0,1,2\}$$

$$\text{Feature} \quad x_j \equiv \begin{matrix} x_j^{0,0} & x_j^{0,1} & x_j^{0,2} \\ x_j^{1,0} & x_j^{1,1} & x_j^{1,2} \\ x_j^{2,0} & x_j^{2,1} & x_j^{2,2} \end{matrix} \quad \{0,1\}$$

$$\{0,1,2\}$$

$$\begin{cases} \log \dfrac{Pr(z_i(1)=1|\mathcal{X}(i),\mathcal{Z}_{\mathcal{N}_i})}{Pr(z_i(0)=1|\mathcal{X}(i),\mathcal{Z}_{\mathcal{N}_i})} = \displaystyle\sum_{j=1}^{m} \sum_{t \in \mathcal{D}_{x_j}} \beta_j^{1,t} x_j^{1,t} + \sum_{k \in \mathcal{N}_i} \sum_{s \in \mathcal{D}_{z_k}} \eta_{1,s} z_k(s) \\[3ex] \log \dfrac{Pr(z_i(2)=1|\mathcal{X}(i),\mathcal{Z}_{\mathcal{N}_i})}{Pr(z_i(0)=1|\mathcal{X}(i),\mathcal{Z}_{\mathcal{N}_i})} = \displaystyle\sum_{j=1}^{m} \sum_{t \in \mathcal{D}_{x_j}} \beta_j^{2,t} x_j^{2,t} + \sum_{k \in \mathcal{N}_i} \sum_{s \in \mathcal{D}_{z_k}} \eta_{2,s} z_k(s) \end{cases}$$

$$Pr(z_i(0)=1) = 1 - Pr(z_i(1)=1) - Pr(z_i(2)=1)$$

**Scalability**

Training Time in sec.(Log) vs Grid Size

- TurboReg
- ngspatial

**Accuracy**

Accuracy [0, 1] vs Grid Size

- TurboReg
- ngspatial

At **least three orders of magnitude** performance gain, while accuracy is **almost the same**.

# Spatial Probabilistic Graphical Modeling (SPGM)

- Performing *uncertain* (i.e., prob.) predictions over spatial data
  - ❑ Classical ML approaches (e.g., regression) ignore the probabilistic relationships

**Disaster Analysis**  **Crime Analysis**  **Public Health Monitoring**  **Geo-tagged Ads**



- Representing the world as a collection of *random variables* with joint probabilistic distribution
  - ❑ Tasks: learning the distribution, and inferring unknown variables via the distribution



**Spatial Markov Random Field (SMRF)**  **Spatial Hidden Markov Model (SHMM)**  **Spatial Bayesian Network (SBN)**

# SMLN for SPGM

■ Generates an equivalent set of weighted rules containing logical predicates for any SPGM input

 ❑ Weights represent the original SPGM parameters

 ❑ Rules follow the syntax of the DDlog logic programming framework

**Spatial Markov Random Field (SMRF)**

**Spatial Hidden Markov Model (SHMM)**

**Spatial Bayesian Network (SBN)**

| MLN Rules |
|---|
| $[P_1 \wedge F_1, \beta_1]$ |
| $[P_1 \wedge P_2, \eta]$ |
| $[P_1 \wedge P_3, \eta]$ |
| $[P_2 \wedge F_2, \beta_1]$ |
| $[P_2 \wedge P_4, \eta]$ |
| ……. |

| MLN Rules |
|---|
| $[O_1 \rightarrow P_1, b]$ |
| $[P_1 \rightarrow P_2, a]$ |
| $[O_2 \rightarrow P_2, b]$ |
| $[P_2 \rightarrow P_3, a]$ |
| $[O_3 \rightarrow P_3, b]$ |
| ……. |

| MLN Rules |
|---|
| $[!P_1 \vee !F_1 \vee !C_1]$ |
| $[!P_3 \vee !P_1 \vee !F_3 \vee !C_1]$ |
| $[!P_2 \vee !F_2 \vee !C_1]$ |
| $[!P_4 \vee !P_2 \vee !F_4 \vee !C_1]$ |
| $[!D_1 \vee !F_1]$ |
| ……. |

# Machine Learning meets Big Spatial Data

# Routing..



**Routing Algorithm**

**Source**

**Destination**

**Route**

Precomputed Routes

Graph D

**Map**

# Routing..



## Routing Algorithm

**Source**

**Destination**

Precomputed Routes

**Path**

**Route**

**Estimated Time of Arrival (ETA)**

**Topology**

**Metadata**

**Map**

UNIVERSITY OF MINNESOTA

# QARTA: An ML-based System for Accurate Map Services

- **Map-Centric:** QARTA *learns* its own map in terms of topology and metadata

- **Query Calibration:** QARTA *learns* the error margins of various algorithms and use it to calibrate its answer

UNIVERSITY OF MINNESOTA

# QARTA: Why..??

- Problem came up from the Taxi company working in Qatar


KARWA TAXI APP

CACM, April 2021

Too much construction and road changes in town (in preparation to FIFA 2022)

Commercial maps cannot cope with such changes in road networks, and are not cheap

### Traffic Routing in the Ever-Changing City of Doha

BY SOFIANE ABBAR, RADE STANOJEVIC, SHADAB MUSTAFA, AND MOHAMED MOKBEL



**The Peninsula** — QATAR'S DAILY NEWSPAPER
Local focus, Global vision
3rd Best News Website in the Middle East in 2017

## Qatar road network increased three times between 2013-18: Ashghal

24 Apr 2018 - 11:58

Al Muhannadi said that the length of the road network increased by about three times between 2013 and 2018 compared to before 2013. He said that the volume of roadworks carried out over the past five years also increased from 1,700 km to 6,000 kilometers, while sanitation capacity doubled, rainwater drainage grew 7 times, and pedestrian trails increased 12 times during the same period.



أشغال — خرائط جوجل

Raya Daily (Sept. 8, 2020), 20

**traffic** TECHNOLOGY TODAY.COM

| Reason | % |
|---|---|
| Maps recommend routes that take longer | ~44 |
| Map shows a location far away from actual drop-off point | ~43 |
| Delivery entrance is not the same place as the street address | ~40 |
| Maps don't know about road closures | ~34 |
| Time estimates are inaccurate | ~28 |
| Updated details are not reflected (new exits, turn lanes, one-way streets, etc.) | ~27 |
| Road names are incorrect | ~16 |
| Other | ~2 |
| My mapping apps are always perfect | ~5 |

% 0 5 10 15 20 25 30 35 40 45

## Poor maps costing delivery companies US$6bn annually 💬 0

BY ADAM FROST ON FEBRUARY 19, 2020                                    MAPPING

**Based on a survey of delivery drivers in the USA and conducted by an independent research firm, the first 'Mapping in Logistics Report' has revealed that 'broken maps' are costing the logistics sector an estimated US$6bn annually.**

# Edge Weight Inference: Who is doing it?

- Traffic departments: Loop detectors or plate recognition



**Distributed Information System**

**Traffic Police Station**

**Traffic Control Server Local Database**

- Co...

**Edge Weights are considered as proprietary information, not to be shared**

with

lyft

UBER

**99 phones and a little red wagon**

*The streets were mostly empty, but the map showed a traffic jam*

By Jay Peters | @jaypeters | Feb 3, 2020, 5:08pm EST

# Edge Weight Inference in QARTA:

- **Input:** Trips (Pickup time/location, Drop off time/location)

(A, F, 15) ➔ $w_2 + w_5 + w_6$ = 15

(B, H, 28) ➔ $w_3 + w_7 + w_8 + w_9 + w_{11}$ = 28

(A, I, 19) ➔ $w_1 + w_3 + w_7 + w_8 + w_9$ = 19

…

- **Objective:** Given a set of edges, each with length $l_e$ and unit *length weight* $W_e$, a set of trips $T$, each with a path $P_t$, find $W_e$ that minimize:

$$\sum_{t \epsilon T} \left( \sum_{e \epsilon Pt} W_e l_e - \delta_t \right)^2$$
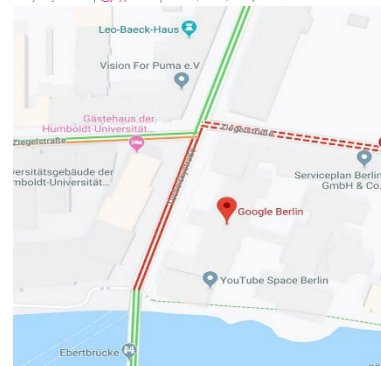
X equations in Y unknowns → **Ridge Regression Analysis** → Edge weights per granularity

Temporal Granularity

- **Challenges:**
  - ❑ A direct solution may result in zero or negative weights
  - ❑ Scalability is a major issue: Hundreds of thousands of edges with millions of trajectories
  - ❑ Over-fitting for unreliable edges
  - ❑ Need to accommodate for a fine granularity (e.g., 168 hours per week)

# Edge Weight Inference in QARTA

- After several tuning steps: (Details in the paper)

- **Objective:**

Regularization strength

Average speed

$$\sum_{t \epsilon T} \left( \sum_{g:P_t \cap Hg \neq \emptyset} W_g L_g + W_0 \sum_{e \, \epsilon \, (Pt \backslash H)} l_e - \delta_t \right)^2 + \alpha \sum_g (W_g - \sigma)^2$$

$$L_g = \sum_{e \, \epsilon \, Hg} l_e$$

X equations in Y unknowns

X` equations in Y` unknowns

Edge weights per granularity

**Tuning Steps**

**Ridge Regression Analysis**

10K equations in 500K unknowns

1K equations in 5K unknowns

Edge weights per hour

Granularity    Hour

# Estimated Time of Arrival (ETA)

- The accuracy of query answers heavily rely on Estimated time of Arrival

**Idea:** Can we study the error patterns of each algorithm under various context, and use to adjust the query answer.


OSRM


Google Maps

## VB
Uber taps ClimaCell to improve ETA estimates with hyper-local weather data

PAUL SAWERS  @PSAWERS  FEBRUARY 6, 2020 8:00 AM

Uber is partnering with weather technology company ClimaCell to enable more accurate estimated time of arrival (ETA) predictions for drivers and riders.

Founded in 2016, Boston-based ClimaCell specializes in real-time weather forecasts. Rather than relying on government data typically garnered from

# Model Building

- Trip: (Pickup time/location, Drop off time/location, $\delta$ )
  - $\delta$ is the difference between actual and estimated time of the trip

Trip Data → **Feature Extraction** → Feature Vector $V$ → **Training Data** ( $V, \delta$ ) → **Gradient Boosting Regressor** → Model $\mathcal{M}$ that maps $V$ to $\delta$

- Features in $V$ that impact $\delta$
  - Spatial Zoning
    - Origin
    - Destination
  - Temporal Zoning
    - Pickup time
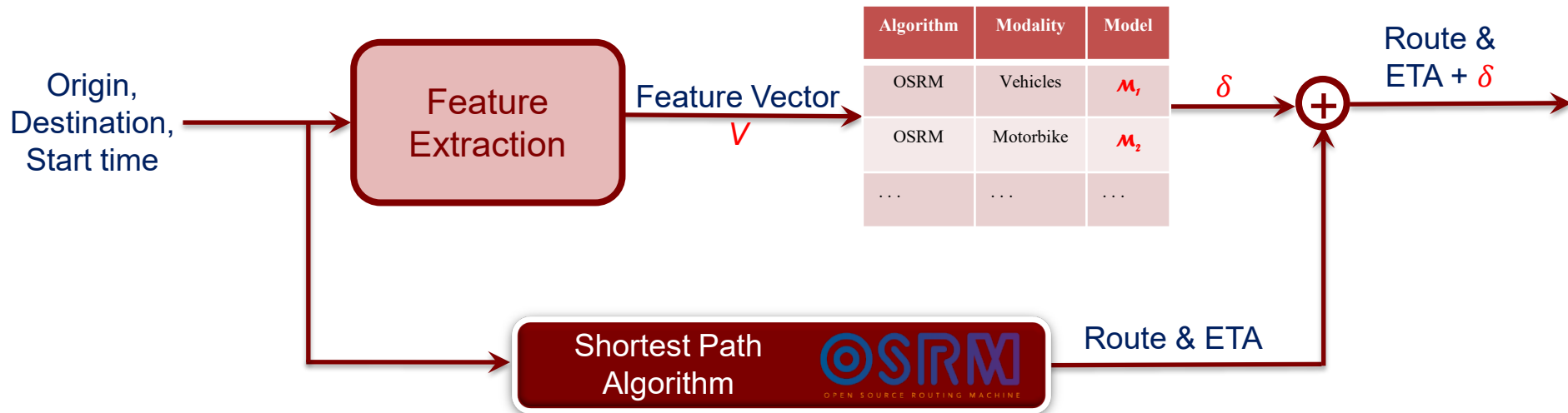    - Drop off time
  - Trip Characteristics
    - Trip distance
    - Trip duration

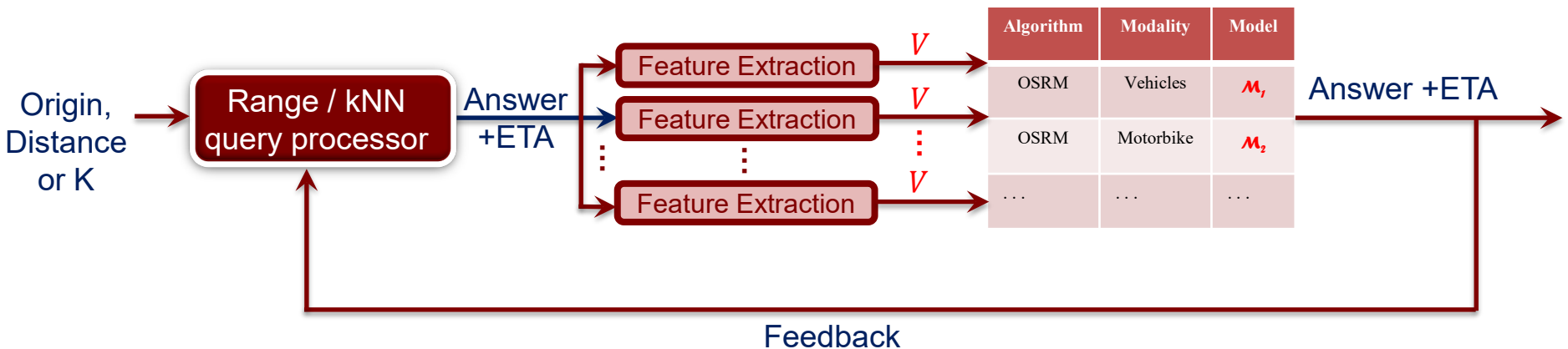- A model $\mathcal{M}$ will be built for each ETA algorithm and driving modality

| Algorithm | Modality | Model |
|-----------|----------|-------|
| OSRM | Vehicles | $\mathcal{M}_1$ |
| OSRM | Motorbikes | $\mathcal{M}_2$ |
| . . . | . . . | . . . |

# Query Calibration in QARTA

■ Shortest Path queries



Origin, Destination, Start time → **Feature Extraction** → Feature Vector $V$

| Algorithm | Modality | Model |
|-----------|----------|-------|
| OSRM | Vehicles | $\mathcal{M}_1$ |
| OSRM | Motorbike | $\mathcal{M}_2$ |
| . . . | . . . | . . . |

$\delta$ → ⊕ → Route & ETA + $\delta$

**Shortest Path Algorithm** OSRM → Route & ETA

■ Range and kNN queries



Origin, Distance or K → **Range / kNN query processor** → Answer +ETA → Feature Extraction $V$ / Feature Extraction $V$ / Feature Extraction $V$

| Algorithm | Modality | Model |
|-----------|----------|-------|
| OSRM | Vehicles | $\mathcal{M}_1$ |
| OSRM | Motorbike | $\mathcal{M}_2$ |
| . . . | . . . | . . . |

Answer +ETA

Feedback

# QARTA in Deployment

QARTA is deployed in *all* Taxis in Qatar $\sim 4K$ vehicles

A local food delivery company $\sim 3K$ motorbiks
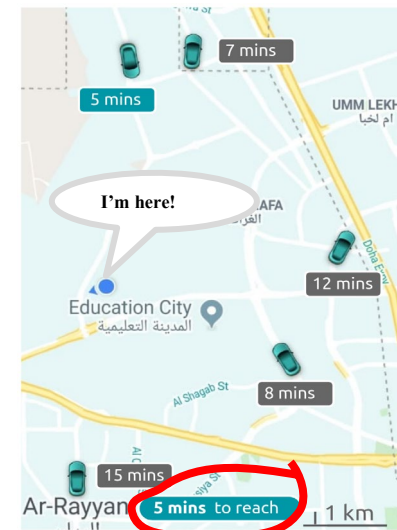
- **QARTA receives:**
  - ❑ $\sim 235K$ daily API calls
  - ❑ $\sim 1 \, Million$ daily GPS tracks

- **APIs & Services:**
  - ❑ In-traffic routes
  - ❑ Travel time estimation
  - ❑ Complex route planning
  - ❑ OD matrices
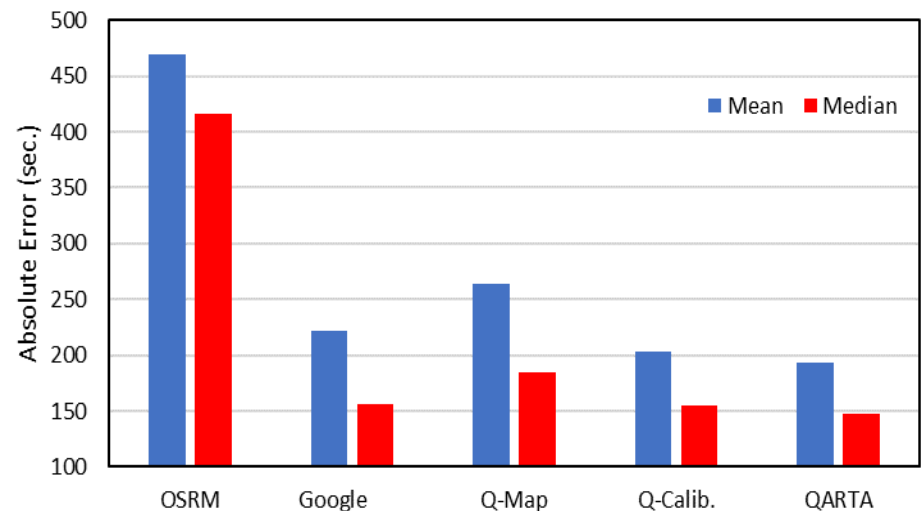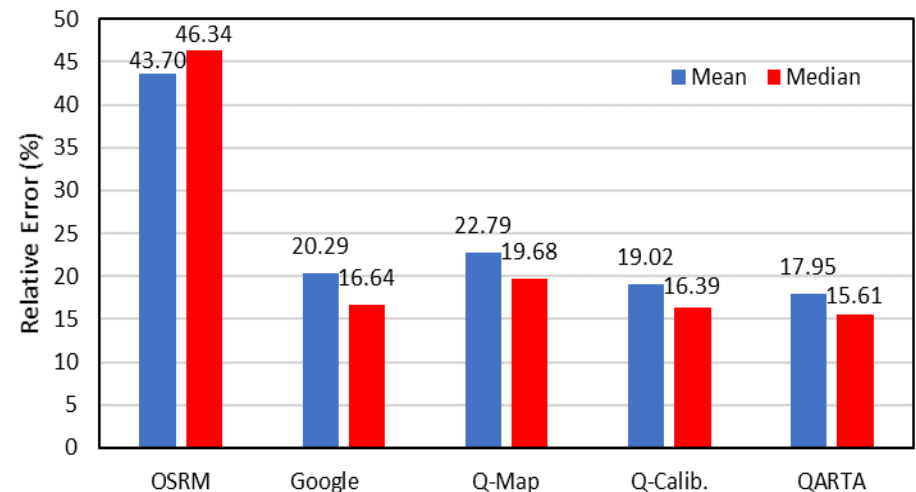  - ❑ Search & addresses

**Routing**



**Fare estimation**

**Taxi Dispatching**

Link: https://qarta.io

# QARTA vs Other Map Services: Shortest Path Query

- **Q-Map**: Runs QARTA Map Making layer without any calibration
  - ❑ OSRM on QARTA map

- **Q-Calib**: Runs QARTA calibration without Map Making layer
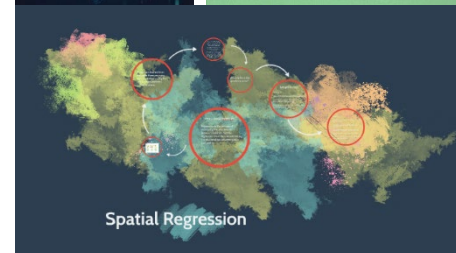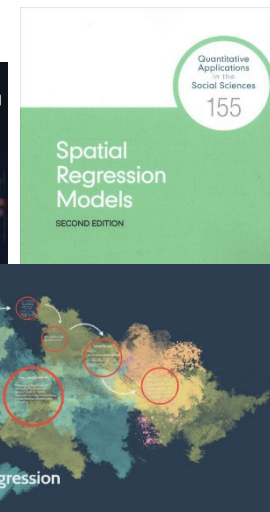  - ❑ Calibrating OSRM engine

# Summary:



**Applications**

**Spatial**

Routing

**Non-Spatial**

CYBER SECURIT

digitalhealth
news • networks • intelligence •

edureka!
TOP 10
APPLICATIONS OF
MACHINE LEARNING

Machine Learning Applications
Retail · Travel · Healthcare · Finance · Media

Knowledge Base

SPATIAL REGRESSION MODELS FOR THE SOCIAL SCIENCES

Quantitative Applications in the Social Sciences 155
Spatial Regression Models
SECOND EDITION

Spatial Regression

Knowledge Base

**Non-Spatial**          **Spatial**

**ML Fundamental Algorithms**

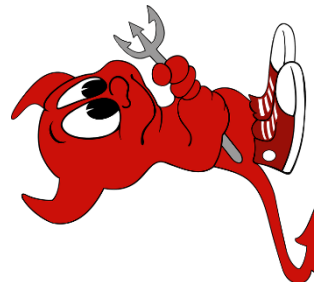# Machine Learning meets Big Spatial Data



ETHICS

*Machine Learning*

*Big Spatial Data*

PRIVACY

POLICIES

Thank you