Visual Exploration of Big Spatio-Temporal Urban Data: A Study of New York City Cab Trips

VAST INFOVI BIOV

Nivan Ferreira, Jorge Poco, <u>Huy T. Vo</u>, Juliana Freire, and Claudio T. Silva Poly and CUSP, New York University

ata: Trips Idio T. Silva ork University

Big Cities – the world is urbanizing



Cities are the cause of our problems and the source of the solutions









Exploring Urban Data

Infrastructure



Condition, operations

Meteorology, pollution,

Environment

noise, flora, fauna



Relationships, location, economic /communications activities, health, nutrition, opinions, ...

Properly acquired, integrated, and analyzed, data can

- Take government beyond imperfect understanding
 - Better (and more efficient) operations, better planning, better policy
- Improve governance and citizen engagement
- Enable the private sector to develop new services for citizens, governments, firms
- Enable a revolution in the social sciences







Urban Data Sources

- Organic data flows

- Administrative records (census, permits, ...)
- Transactions (sales, communications, ...)
- Operational (traffic, transit, utilities, health system, ...)
- New and social media (Twitter feeds, blog posts, Facebook, ...)

Sensors

- Personal (location, activity, physiological)
- Fixed *in situ* sensors
- Crowd sourcing (mobile phones, ...)
- Choke points (people, vehicles)
- Opportunities for "novel" sensor technologies
 - Visible, infrared and spectral imagery
 - RADAR, LIDAR
 - Gravity and magnetic, seismic, acoustic
 - Ionizing radiation, biological, chemical











Taxis as Sensors for NYC

Taxis are sensors that can provide unprecedented insight into city life: economic activity, human behavior, mobility patterns, ...

"What is the average trip time from Midtown to the airports during weekdays?" "How the taxi fleet activity varies during weekdays?"

"How was the taxi activity in Midtown affected during a presidential visit?"

"How did the movement patterns change during Sandy?"

"Where are the popular night spots?"







Exploring Taxi Trips: Challenges

Taxi data are:

- **Big:** 520 million trips -- ~500k trips/day
 - Can't use existing tools for interactive exploration
- **Complex**: Multiple variables: *spatial and temporal + trip attributes*
 - Hard to select data -- too many data slices
- **Dirty:** Taxis in the river...

Domain scientists and decision makers are unable to interactively explore the whole data







Exploring Taxi Trips: Challenges

Taxi data are:

- **Big:** 520 million trips -- ~500k trips/day
 - Can't use existing tools for interactive exploration
- **Complex**: Multiple variables: *spatial and temporal + trip attributes*
 - Hard to select data -- too many data slices
- **Dirty:** Taxis in the river...





TaxiVis: Visually Exploring NYC Taxi Data

- New model that allows users to visually query taxi trips, easily select and compare different spatial-temporal slices
 - Data selection through visual manipulations
 - Use visualization to explore selected data
- Support for origin-destination queries that enable the study of mobility across the city
- Use multiple coordinated views to allow comparisons, and brushing to support query refinements
- Use of adaptive level-of-detail rendering and heat maps to generate clutter-free visualization for large results
- Scalable system that provides interactive response times for spatio-temporal queries over large data





Related Work – taxi data

- Recommendation Systems [*Ge et al. 2010*] [*Yuan et al. 2011*]
- Land-use Classification
 [Pan et al. 2013]
- Human Mobility [Veloso et al. 2011] [Liang et al. 2012] [Peng et al. 2012]













Related Work – mobility visualization

Flow Maps [Phan et al. 2005]

OD Maps [Wood et al. 2010]











Related Work – querying and visualizing spatio-temporal data

- Spatio-temporal querying model [Peuquet 1994]
- Cross-filtered views [Weaver 2008]
- Visual query languages
 - GeoPQL [Ferri and Rafanelli 2005]
 - Moving GeoPQL [D'Ulizia et al. 2012]
 - Query-by-trace [Erwig and Schneider 2000]
 - TrajectoryLenses [Krueger et al. 2013]

Dyadic Event Data Sets APRE 4007 exet in the (1000) RELE 4007 exet in the (1000) APRE 4007 exet in (19239) APRE 4007				Time Series (All Event			ent	
		Events By Code				Jern		
Code 1	Weight	Times	O C					_
31	1.0	· · ···· · · · · · · · · · · · · · · ·	422 *	0.1,21.	er an de set		S. Andrewski	ŝ.,
33	2.8	· · · · · · · RECEIVE · · ·	309	1.1	1000	204	- 1. Artis	
223	-10.0		245			-		
212	.9.0	ARREST PERSON	182					-
23	.0.2	NEUTRAL COMMENT	169	303.01.29	2003.02.18	20	03.03.10 2003.0	13.30
20	-0.6		100	to a				
32	1.9	VISIT	156	POL	BAR >		5	~
IRQBUS IRQCHR IRQELI IRQGOV		india de la construction de la cons La construction de la construction d	3394 6 3 145 836 12	LBY	Er	SAU	IRN PAK	>
🗌 Filte	r Codes	Filter Targets	lter Time	<u></u>	EI	15		
		Events By Targets		Filter 0	lodes 🖌	Filter S	iources Filte	r Ta
Val	ue	Times	0 C				Events	
USAMIL	(62 *	Date S	lource Target	Code	Tag	
USA	(31	2003.03.30	N IRQ	70	(Provide aid)	
MOSSHI	ĺ		22	2003 03 30 IR 2003 03 30 IR	AUS USAMIL	160	(Make a visit) (Reduce relations)	-
JOR	Ì		11	2003.03.30 U	IRQ IRQ	73	(Provide humanitarian aid)	-
USAGOV	, ì		10	2003-03.30 0	USAMIL	23 223	(Host a visit)	-
190			10	2003.03.30 P	SE ISRCOP	23	(Host a visit)	
				2003 03 30 N	JOANN GBRGOV	30	(Approve)	\neg
				2003 03 30 P	10070 (0044)	223	0	-
2 Filte	e Codes	2 Filter Sources	tec Time	-	PLOED TOPPET		0	







When

Desiderata in TaxiVis

End-to-end solution for interactive visual analytics

- Couples database back-end with a usable interface front-end
- Support interactive queries for the entire (growing) datasets
 - Out-of-core data access







Data – limits of existing technology

Raw data:

- 3 years: 2009, 2011, and 2012
- 150 GB in 48 CSV files
- 520M trips total
- After cleanup and transformation:
 - 50GB in binary format
 - 12 fields with 2 temporal spatial attributes

	SQLite	F
Storage Space in GB	100	
Building Indices in Minutes (One Year of Data)	3,120	
1K Items Query in Seconds	8	
100K Items Query in Seconds	85	



PostgreSQL	
200	
780	
3	

24



Data – limits of existing technology

	Raw data:		SQL ite	[
	• 3 years: 2009, 2011, and 2012			-
	 150 GB in 48 CSV files 	Storage Space in GB	100	
•	 520M trips total After cleanup and transformation: 	Building Indices in Minutes (One Year of Data)	3,120	
	 50GB in binary format 12 fields with 2 temporal spatial 	1K Items Query in Seconds	8	
	attributes	100K Items Query in Seconds	85	
•	Solution:			<u> </u>

- New spatio-temporal index based on out-of-core Kd-tree
 - Can also index other attributes



PostgreSQL	Our Solution
200	30
780	28
3	0.2
24	2



Data Exploration: A Two-Phase Process



We unify the two phases of the process through visual operations





Interactive Visual Exploration

- Help users easily select data slices
 - Composition of spatio-temporal constraints
- Provide visualizations of the slices within the spatio-temporal context
- Multiple coordinated views
 - Time series, histogram plots
 - Heat maps
- Comparative visualizations
 - Multiple query results
 - Exploration in time
 - Summary of attributes









Visual Query Model

- Data selection by visual operations
- Each data slice can be assigned a different visual representation
 - Spatial context is maintained in the map view
- Query Expressiveness [Peuquet 1994]
 - when + where \rightarrow what
 - when + what \rightarrow where
 - where + what \rightarrow when











When + Where \rightarrow What







When + Where -> What





When + Where \rightarrow What





When + Where \rightarrow What









TaxiVis in Action







Cross-filtered Views: trip distributions by area







Temporal Comparison: Hurricane Sandy









Analyzing Movement







Comparative Visualization: Night Life in NYC





Saturday

Monday

1

ø _ap



Conclusion & Limitations

Easy-to-use system to interactively explore large multivariate spatial-temporal data

- Couples database and visualization
- **Desktop-based application**
 - Web-based map
- Need support for multiple data layers





Future and Ongoing Work

- Apply to other urban mobility data, e.g., data from the NYC bike share program
- Support additional data layers: weather, gas prices, news, tweets, etc.
- Utilize parallel processing







e program c.



Acknowledgements

- The Taxi & Limousine Commission of New York City
- National Science Foundation (CNS-1229185, CNS-1153503, IIS 1139832, IIS-0905385, IIS-1142013, AGS 0835821)
- The Department of Energy







Thanks!





