

RUTGERS THE STATE UNIVERSITY OF NEW JERSEY

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 26

## MINING DATA FROM MOBILE DEVICES

Applications: Location, Ads, Privacy

Spiros Papadimitriou, *Tina Eliassi-Rad*  
[spiros.papadimtriou.tina.eliassi@rutgers.edu](mailto:spiros.papadimtriou.tina.eliassi@rutgers.edu)



RUTGERS

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 27

### Mobile Real-Time Bidding (RTB) Ad Ecosystem

Advertisers

Trading Desks

Demand-Side Platforms

Ad Exchanges

Supply-Side Platforms

Data Partners

Dynamic Creative Optimization Partners

Ad Verification and Brand Protection

PUBLISHERS

CONSUMERS

businessinsider.com

RUTGERS

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 28

### Ad Targeting on Mobile Devices

- Need reliable location information
  - We collected 21.6M RTB requests on Wed 2/6/2013
  - The majority of them (57%) did **not** have location information
- Reasons behind the missing location information
  1. The RTB system did not forward it
  2. The SSP did not forward it
  3. The device did not capture it
  4. The user did not enable location-based services

- RTB: Real-Time Bid
- SSP: Supply Side Provider

RUTGERS

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 29

### Problem Definition

- **Accurately infer location information for IPs in RTB requests**
- What type of location information should we infer?
  - Latitude, longitude
  - Census Block Groups (CBGs) ← *privacy-friendly location information*
  - Zip codes
  - ...

RUTGERS

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 30

### Census Block Groups (CBGs)

- Assumption
  - CBGs comprise location information fine-grained enough for useful hyper-local ad targeting, yet coarse-grained enough to avoid major privacy concerns.
- Why is this reasonable?
  - Covers a contiguous area
  - Never crosses state or county boundaries
  - Contains between 600 and 3,000 people
  - US is divided into ~212K CBGs

- US is divided into ~8.2M CBs
- US is divided into 43K zip codes

RUTGERS

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 31

### Geo-locating IP addresses on mobile networks is hard

- Balakrishnan et al. (IMC 2009) examined properties of cell-phone IP addresses
- Mobile IPs are ephemeral and their addresses are itinerant
- Example: An individual cell phone can report different IP addresses to various servers within a short time-period



Answers to IP → Location queries provided by 7 geo-location services; the actual cell phone is in Mountain View, CA.

M. Balakrishnan, I. Mohomed, and V. Ramasubramanian, Where's that phone? Geolocating IP addresses on 3G networks. In IMC, pages 294–300, 2009.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 32

## Many devices often use the same public IP address

- Metwally & Paduano (KDD 2011) estimated the number of users of an IP address by keeping track of the application-specific traffic
- The primary goal of their work was to combat "abusive" traffic (such as DDoS attacks, ad click fraud and email spam) without violating the user privacy

Frequency  
 10<sup>0</sup> 10<sup>1</sup> 10<sup>2</sup> 10<sup>3</sup> 10<sup>4</sup> 10<sup>5</sup> 10<sup>6</sup>  
 1 10 100 1000 10000  
 Est. Number of Distinct Count. (User-IPs)  
 Estimated Size of 10M Random IPs

A. Metwally and M. Paduano. Estimating the number of users behind IP addresses for combating abusive traffic. In KDD, pages 249–257, 2011.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 33

## Hyperlocal: A Graph Mining Solution

- Classify IPs as mobile vs. non-mobile
- Construct a movement graph with mobile and non-mobile IP nodes
- Use a local relational classifier on each unknown node to infer latitude and longitude
- Assign Census Block Group (CBG) ID to the inferred latitude and longitude using a  $k$ -nearest neighbor approach

- L.T. Le, T. Eliassi-Rad, F. Provost, L. Moeres. *Hyperlocal: Inferring Location of IP Addresses in Real-time Bid Requests for Mobile Ads*. ACM SIGSPATIAL LBSN 2013.
- F. Provost, T. Eliassi-Rad, L. Moeres. *Methods, Systems, and Media for Determining Location Information from Real-time Bid Requests*. US Patent Number 9014717, issued April 21, 2015.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 34

## 1. Classifying IPs as Mobile vs. Non-mobile

```

graph TD
    IP --> Q1{Does 3rd party database say IP is mobile?}
    Q1 -- Yes --> M1[IP is classified as mobile]
    Q1 -- No --> Q2{Does the IP appear outside radius r in 24 hours?}
    Q2 -- Yes --> M2[IP is classified as mobile]
    Q2 -- No --> NM[IP is classified as non-mobile]
  
```

- Mobile IPs tend to change position more quickly than non-mobile IPs
- Default  $r = 100m$  (range of current Wi-Fi routers)

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 35

## Reversing Third Party Decision

Radius (in km)	% IPs Reversed from 3rd Party Label (Oct-2012)	% IPs Reversed from 3rd Party Label (Feb-2013)
0	~15%	~20%
0.05	~11%	~13%
0.1	~11%	~13%
0.5	~10%	~12%
1	~9%	~11%
10	~5%	~5%
50	~3%	~3%
Infiniti	~1%	~1%

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 36

## 2. Constructing IP $\times$ IP Movement Graph

- Circle nodes are non-mobile IPs
- Diamond nodes are mobile IPs
  - Mobile IPs are time-stamped because they are transient
- An edge indicates an NUID's movement from one IP to another
  - Each edge stores
    - # of movements between its endpoints, and
    - inter-arrival times (IATs) for all movements across it

NUID: A network exchange ID associated with each device

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 37

## 3. Employing Local Relational Classifiers

- Why local?
  - The farther out one moves in the movement graph, the farther away one gets geographically
  - The movement graph is big  $\rightarrow$  non-local approaches can be computationally burdensome
- Local relational classifier used:  $wvRN$ 
  - $wvRN$  stands for weighted-vote Relational Neighbor [Macskassy & Provost, JMLR 2007]
  - $wvRN$  estimates the class membership probabilities and assumes homophily in the network data

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 38

## What weights should we use in $wvRN$ ?

- Number of movements
  - Intuition: A node  $v$  will be closer in distance to its neighbors with whom it has more movements
  - The data contains many edges with only one movement
- Minimum IAT
  - Intuition: The longer the IAT, the longer distance the user has potentially moved (sans traffic)
  - Uses the normalized minimum IAT between two IPs
  - Weight on the movement between node  $v$  and its neighbor  $i$  is
 
$$w_i = \frac{\min IAT_v}{\min(t, \forall t \in IAT(v, i))}$$
  - $\min IAT_v$  is the minimum IAT across all of a given node  $v$ 's edge

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 39

## $wvRN$ Equations

- Inputs
  - Node  $v$
  - Its neighbors  $Nbr(v)$
  - Weights  $W$  on the edges between  $v$  and its neighbors
- $wvRN$  equations
 
$$latitude(v) = \frac{\sum_{i \in Nbr(v)} w_i \times latitude(i)}{\sum_{i \in Nbr(v)} w_i}$$

$$longitude(v) = \frac{\sum_{i \in Nbr(v)} w_i \times longitude(i)}{\sum_{i \in Nbr(v)} w_i}$$

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 40

## 4. Assigning CBGs as Proxies for Location

- Infer location of a hashed public IP address at the CBG level and not at the  $\langle$ latitude, longitude $\rangle$  level
- Why use CBG ?
  - It provides a more consistent labeling (as in location) for IPs
  - It allows incorporation of external data that uses census data such as demographics
  - In the majority of mobile applications, this level of location information is sufficient for a successful campaign

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 41

## A $k$ -Nearest Neighbor Approach for Assigning a CBG ID to a $\langle$ lat, lon $\rangle$

<ul style="list-style-type: none"> <li>Inputs           <ul style="list-style-type: none"> <li>Location of interest               <ul style="list-style-type: none"> <li><math>loc = \langle lat, lon \rangle</math></li> </ul> </li> <li>For each CBG <math>i</math> in the US,               <ul style="list-style-type: none"> <li><math>i</math>'s centroid: <math>c_i = \langle lat_i, lon_i \rangle</math></li> <li><math>i</math>'s area <math>a_i</math> in km</li> </ul> </li> </ul> </li> <li>Output           <ul style="list-style-type: none"> <li>The CBG ID that contains <math>loc</math></li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Procedure           <ol style="list-style-type: none"> <li><math>C =</math> centroids of the <math>k</math> nearest CBGs to <math>loc</math></li> <li>For <math>j</math> in <math>C</math> <ul style="list-style-type: none"> <li>Calculate the distance between <math>loc</math> &amp; centroid of the <math>j^{th}</math> nearest CBG               <ul style="list-style-type: none"> <li><math>d_j = distance(loc, c_j)</math></li> </ul> </li> <li>Calculate the radius of CBG corresponding to the <math>j^{th}</math> centroid               <ul style="list-style-type: none"> <li><math>r_j = \sqrt{a_j / \pi}</math></li> </ul> </li> <li>Calculate the ratio of distance over radius               <ul style="list-style-type: none"> <li><math>ratio_j = d_j / r_j</math></li> </ul> </li> </ul> </li> <li>Return the CBG ID corresponding to <math>\min(ratio_j)</math>, for all <math>j</math> in <math>C</math></li> </ol> </li> </ul>
--	--

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 42

## Recap of *Hyperlocal*\*

- Classify IPs as mobile vs. non-mobile
- Construct a movement graph with mobile and non-mobile IP nodes
- Use a local relational classifier on each unknown node to infer latitude and longitude
- Assign Census Block Group (CBG) ID to the inferred latitude and longitude using a  $k$ -nearest neighbor approach

\* <http://eliassi.org/papers/le-ibsn13.pdf>

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 43

## Experiments

- Experiments are divided into nine combinations of *infer location for X using Y*
- Values for  $X$  are 'all IPs', 'mobile IPs', and 'non-mobile IPs'
- Values for  $Y$  are 'all neighbors', 'mobile neighbors', and 'non-mobile neighbors'
- Measure **accuracy** by checking the predicted CBG ID vs. the actual CBG ID of an IP

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 44

## Implementation & Runtime

- Hardware & OS: Macbook Pro with
  - CPU 2.66 GHz Intel Core i7
  - RAM 8 GB DDR3
  - hard drive 500 GB SSD
  - OS X 10.8
- Language: Python
- Supporting Software: NetworkX & MongoDB
- Runtime: On average 1.2 milliseconds to process each RTB request

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 45

## Data

Data Name	Collection Date	# RTB Requests with Valid US NUIDs	% RTB Requests without Location	% RTB Requests from Mobile IPs
Oct-2012	Mon 10/01/2012	44.1M	36.5%	57.3%
Feb-2013	Wed 02/06/2013	21.6M	56.7%	47.7%

- From Oct-2012 to Feb-2013
  - # of RTB requests decreased by ~50%
    - Due to reductions from SSPs
  - # of requests without location information increased by ~55%
  - # of requests from mobile IPs decreased by ~17%

All IPs are hashed public IP addresses.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 46

## IPs on the US Map

Oct-2012

Feb-2013

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 47

## Data Characteristics per SSP

SSP4 provides about 50% of the requests for both datasets.

Supply Side Provider	Oct-2012 (%)	Feb-2013 (%)
SSP1	~5	~5
SSP2	~0	~0
SSP3	~5	~5
SSP4	~50	~50
SSP5	~10	~10
SSP6	~10	~10
SSP7	~30	~10

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 48

## % Requests with and without Location Information per SSP

Oct-2012

Feb-2013

- None of the requests from SSP1 and SSP3 have location information.
- All the requests from SSP7 have location information.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 49

## Oct-2012 Homophily per SSP

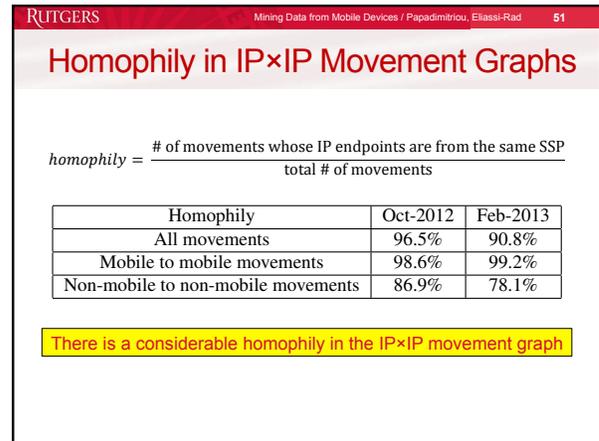
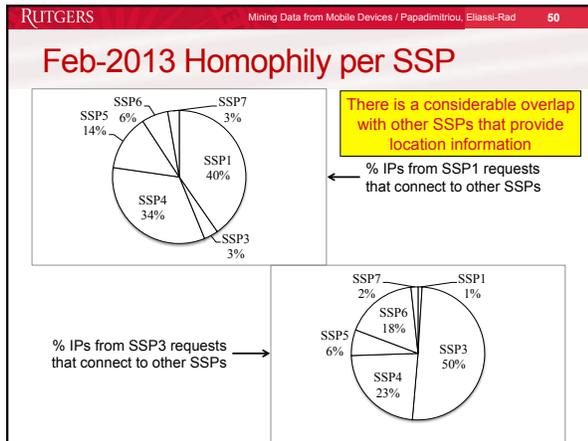
SSP6 1%  
SSP7 21%  
SSP1 40%  
SSP4 34%  
SSP3 2%  
SSP5 2%

There is a considerable overlap with other SSPs that provide location information

← % IPs from SSP1 requests that connect to other SSPs

% IPs from SSP3 requests that connect to other SSPs →

SSP6 6%  
SSP7 12%  
SSP1 1%  
SSP5 3%  
SSP4 27%  
SSP3 51%



RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 52

## Core Results: wvRN(minIAT) vs. wvRN(numMoves)

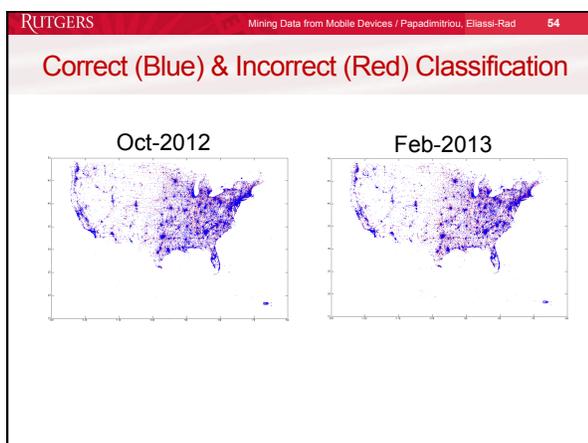
Infer location for all IPs...	Accuracy wvRN(minIAT)		Accuracy wvRN(numMoves)		Number of Predictions		% Mobile in Predictions	
	Oct-2012	Feb-2013	Oct-2012	Feb-2013	Oct-2012	Feb-2013	Oct-2012	Feb-2013
using all neighbors	71.6	70.8	69.7	69.3	1,189,679	1,328,015	91.0%	83.8%
using mobile neighbors	<b>74.4</b>	<b>74.8</b>	<b>72.6</b>	<b>73.5</b>	1,077,644	1,278,674	93.2%	87.4%
using non-mobile neighbors	51.9	57.7	51.7	57.5	98,338	30,777	68.7%	62.1%
Infer location for mobile IPs...								
using all neighbors	74.2	75.0	72.3	73.6	1,082,566	274,900	100%	100%
using mobile neighbors	<b>76.6</b>	<b>79.1</b>	<b>74.7</b>	<b>77.9</b>	1,004,601	243,630	100%	100%
using non-mobile neighbors	50.4	54.7	50.2	54.5	67,553	25,323	100%	100%
Infer location for non-mobile IPs...								
using all neighbors	45.5	49.0	44.0	46.8	107,113	53,115	0%	0%
using mobile neighbors	44.0	45.1	42.5	43.0	73,043	35,044	0%	0%
using non-mobile neighbors	<b>55.4</b>	<b>62.7</b>	<b>55.0</b>	<b>62.0</b>	30,785	15,454	0%	0%

- Differences between the two methods are not statistically significant at the 0.05 level
- Number of predictions varies depending on the particular inference and the neighbor types used in the inference process

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 53

## IATs on movement edges are correlated with distances

- Shorter IAT, shorter distance
- For IPs with only one known neighbor, restricting IAT to  $\leq 60$  minutes
  - Improves accuracy by an average of 12% on Oct-2012 and 23% on Feb-2013 data
- Reduces the number of predictions by an average of 4 times for Oct-2012 and 5 times for Feb-2012
- Restricting IATs to  $> 60$  minutes decrease accuracy



RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 55

## Limitations of a Graph Mining Approach

- Cannot infer location for IPs with no neighbors
  - Use other info – e.g., site visits; subnet info, etc.
- Cannot infer location for IPs with no known neighbors
  - Use collective classification.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 56

### Related Work [Wong et al. NSDI 2007]

- Locate IPs by
  - representing node positions through regions,
  - expressing constraints as areas, and
  - computing locations by solving a system of geometric constraints
- Relies on pings to estimate the round-trip time between two IPs

B. Wong, I. Stoyanov, and E. G. Sirer. Octant: A comprehensive framework for the geolocation of internet hosts. In NSDI, pages 23–23, 2007.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 57

### Related Work [Wang et al. NSDI 2011]

- A client-independent geo-location system
- Like [Wong et al. NSDI 2007] relies on pings to estimate the round-trip time between two IPs
- Also relies on landmarks, which are collected manually

Y. Wang, D. Burgener, M. Flores, A. Kuzmanovic, and C. Huang. Towards street-level client-independent IP geolocation. In NSDI, pages 27–27, 2011.

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 58

### Recap & Open Problems

- Graph mining on just the structure of an IP×IP movement graph to *infer locations*, in terms of CBGs, for *hashed public IP addresses* produces an *accuracy of ~75%*
- Results are impressive since estimating the correct CBG is out of 212K possibilities
- **Open problems**
  - Inference on truncated IP addresses
  - Constrained collective classification

RUTGERS THE STATE UNIVERSITY OF NEW JERSEY Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 59

## MINING DATA FROM MOBILE DEVICES

Applications: Location, Ads, Privacy

Spiros Papadimitriou, *Tina Eliassi-Rad*  
[spiros.papadimitriou, tina.eliassi@rutgers.edu](mailto:spiros.papadimitriou, tina.eliassi@rutgers.edu)

RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 60

### License

These slides are made available under a Creative Commons Attribution-ShareAlike license (CC BY-SA 3.0):  
<http://creativecommons.org/licenses/by-sa/3.0/>

You can share and remix this work, provided that you keep the attribution to the original authors intact, and that, if you alter, transform, or build upon this work, you may distribute the resulting work only under the same or similar license to this one.

© 2013, 2015 Spiros Papadimitriou, Tina Eliassi-Rad