

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 22

MINING DATA FROM MOBILE DEVICES

Applications: Location, Ads, Privacy

Spiros Papadimitriou, *Tina Eliassi-Rad*
spiros.papadimitriou.tina.eliasirad@rutgers.edu

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 23

Mobile Real-Time Bidding (RTB) Ad Ecosystem

The diagram illustrates the Mobile Real-Time Bidding (RTB) Ad Ecosystem. It is divided into several categories:

- Advertisers:** Includes Trading Desks (e.g., CADREON, VARIQMEDIA, b3, IProspect, VIVAMI, IMPACT, KAXIS, TACLIEN, DRAFTFCB), Demand-Side Platforms (e.g., adform, lucidmedia, MediaMath, DIGIANT, admobi, TRIGGIT, rocketHub, Strikeart, DataXu, appnexus, thetrade desk, WDA, TURN, JumpTag, Brightroll, adfonic), Ad Exchanges (e.g., Doubleclick, Microsoft Advertising Exchange, RIGHTMEDIA, FLURRY, Rubicon, pulsepoint, INDECN, smaato, addcloud, zENYXA, adbroker), and Mobile Video RTB (e.g., BrightRoll, YuMea).
- Supply-Side Platforms:** Includes Admeld, PubMatic, rubicon, IMPROVE DIGITAL, ADTECH, ADMETA, AdMarvel, mpub, 24/7, and AdIQity.
- Other Intermediaries:** Data Partners (e.g., action, Almond, bize, Experian, quantcast, LOTAME, DEUSCAR, Viz), Dynamic Creative Optimization Partners (e.g., adacado, adready), and Ad Verification and Brand Protection (e.g., adsafe, DoubleVerify, admistry, EVIDON).
- Verticals:** Advertisers on the left and Publishers/Consumers on the right.

businessinsider.com

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 24

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

Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 25

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)
 - Zip codes
 - ...

← privacy-friendly location information

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

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

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

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. Mohamed, 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 28

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

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 29

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

* <http://eliassi.org/ESM2013TR.pdf>

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

1. Classifying IPs as Mobile vs. 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 31

Reversing Third Party Decision

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

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 33

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 34

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

$$W_i = \frac{\min IAT_v}{\min(i, 812 IAT(v, i))}$$
 - $\min IAT_v$ is the minimum IAT across all of a given node v 's edge

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

$wvRN$ Equations

- Inputs
 - Node v
 - Its neighbors $Nbr(v)$
 - Weights W on the edges between v and its neighbors
- $wvRN$ equations

$$latitude(v) = \frac{\sum_{i \in Nbr(v)} W_i \square latitude(i)}{\sum_{i \in Nbr(v)} W_i}$$

$$longitude(v) = \frac{\sum_{i \in Nbr(v)} W_i \square longitude(i)}{\sum_{i \in Nbr(v)} W_i}$$

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

4. Assigning CBGs as Proxies for Location

- Infer location of a hashed public IP address at the CBG level and not at the (latitude, longitude) 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 37

A k -Nearest Neighbor Approach for Assigning a CBG ID to a (lat, lon)

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

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

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/ESM2013TR.pdf>

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

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 40

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 41

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 42

IPs on the US Map

Oct-2012

Feb-2013

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

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%	~15%
SSP6	~5%	~10%
SSP7	~30%	~10%

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

% 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 45

Oct-2012 Homophily per SSP

SSP1: 40%

SSP4: 34%

SSP7: 21%

SSP3: 2%

SSP5: 2%

SSP6: 1%

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

% IPs from SSP3 requests that connect to other SSPs

SSP3: 51%

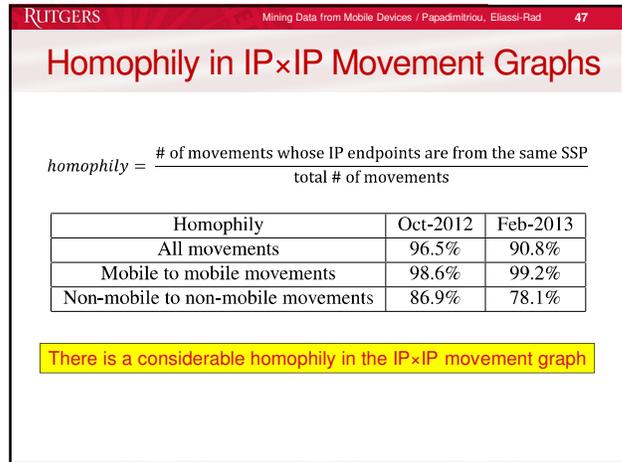
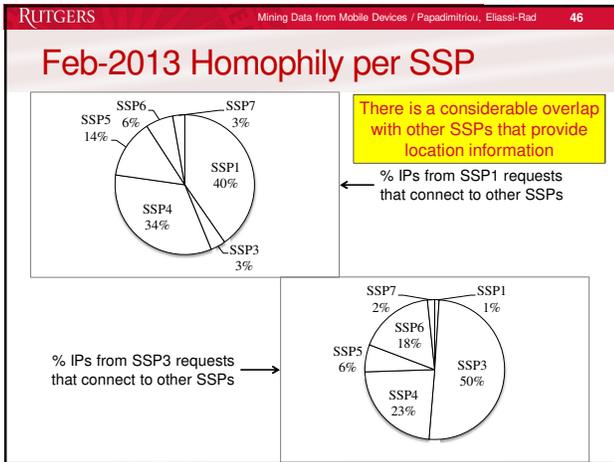
SSP4: 27%

SSP7: 12%

SSP6: 6%

SSP5: 3%

SSP1: 1%



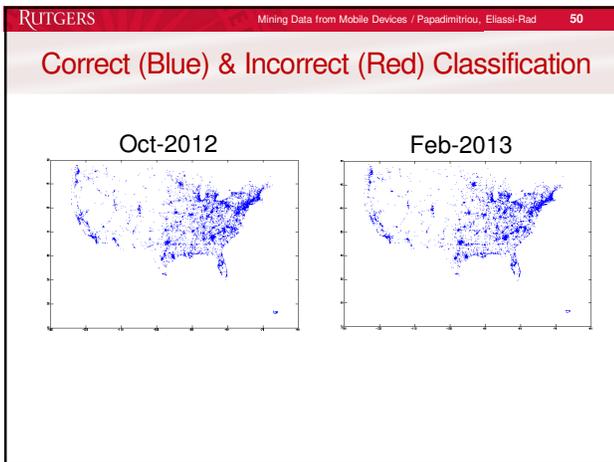
RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 48

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	328,015	91.0%	83.8%
using mobile neighbors	74.4	74.8	72.6	73.5	1,077,644	278,674	93.2%	87.4%
using non-mobile neighbors	51.9	57.7	51.7	55.3	98,338	140,797	68.7%	63.1%
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	76.6	79.1	74.7	77.9	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	55.4	62.7	55.0	62.0	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 49
- ### IATs on movement edges are correlated with distances
- Shorter IAT, shorter distance
 - For IPs with only one known neighbor, restricting IAT to ≤ 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

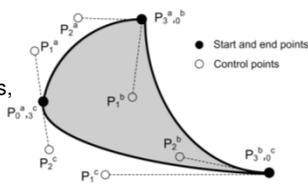


- RUTGERS Mining Data from Mobile Devices / Papadimitriou, Eliassi-Rad 51
- ### Limitations of a Graph Mining Approach
- Cannot infer location for IPs with no neighbors
 - Use other info – e.g., site visits; subnet info, etc.
 - Cannot infer location for IPs with no known neighbors
 - Use collective classification.

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

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

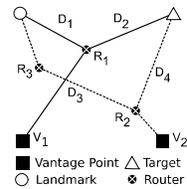


B. Wong, I. Stoyanov, and E. G. Sirer. Octant: A comprehensive framework for the geolocalization of internet hosts. In NSDI, pages 23–23, 2007.

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

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 54

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 55

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 56

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 Spiros Papadimitriou, Tina Eliassi-Rad