

Harvesting Multiple Sources for User Profile Learning: a Big Data Study

Aleksandr Farseev,

Liqiang Nie, Mohammad Akbari, **and Tat-Seng Chua**



What is user profile?



What is human mobility?

• Mobility - contemporary paradigm, which explores various types of people movement.

What is human mobility?

• Mobility - contemporary paradigm, which explores various types of people movement.

- The movement of people
- The quality or state of being mobile



- (Physiology) the ability to move physically
- (Sociology) movement within or between social classes and occupations
- (Chess) the ability of a chess piece to move around the board

Why human mobility?

- Urban planning: understand the city and optimize services
- Mobile applications and recommendations: study the user and offer services





Møbeikitgatatodessevibe pæ?ple





User profile: Mobility + Demography



Multiple sources describe user from multiple views

More than 50% of onlineactive adults use more than one social network in their daily life*

*According Paw Research Internet Project's Social Media Update 2013 (www.pewinternet.org/fact-sheets/social-networking-fact-sheet/)

Multiple sources describe user from multiple views

Research Problems

Multi-source user profiling:

- Geographical user mobility profiling
- User demographic profiling
- Data incompleteness
- Multi–source multi–modal data integration

Multi-source dataset: NUS-MSS*

*http://lms.comp.nus.edu.sg/ research/NUS-MULTISOURCE.htm

NUS-MSS: Data sources

BIGGEST ENGLISH-SPEAKING MICROBLOG BIGGEST PHOTO SHARING SERVICE

NUS-MSS: Data collection

NUS-MSS: Dataset Description

366,268 CHECK-INS

NUS-MSS: Dataset Description

127,276 CHECK-INS

NUS-MSS: Dataset Description

PVV/

304,493 CHECK-INS

NUS-MSS: Dataset Statistics in Singapore

Demographic profiling

User profile: Mobility + Demography

- Linguistic features
 - LIWC
 - User Topics
- Heuristic features
 - Writing behavior

An efficient and effective method for studying the various emotional, cognitive, structural, and process components present in individuals' verbal and written speech samples. Can be highly related to one's demography.

A text analysis software.

Dictionary

Percentage (%)

Word category

- Linguistic features
 - LIWC
 - User Topics
- Behavioral features
 - Writing behavior

Users of similar gender and age may talk about similar topics e.g. female users – about shopping, male – about cars; youth – about school while elderly – about health.

LDA word distribution over 50 topics for collected Twitter timeline.

- Linguistic features
 - LIWC
 - User Topics
- Heuristic features
 - Writing behavior

As we mention from our research – user's writing behavioral patterns are highly correlated with e.g. age (individuals from 10 – 20 years old are making two times less grammatical errors than 20 -30 years old individuals)

Feature name	Description		
Number of hash tags	Number of hash tags mentioned in message		
Number of slang words	Number of slang words one use in his tweets. We calculate number of slang words / tweet and compute average slang usage		
Number of URLs	Number of URL's one usually use in his/her tweets		
Number of user mentions	Number of user mentions – may represent one's social activity		
Number of repeated chars	Number of repeated characters in one tweets (e.g. noooooooo, wahhhhhhh)		
Number of emotion words	Number of words that are marked with not – neutral emotion score in Sentiment WordNet		
Number of emoticons	Number of common emoticons from Wikipedia article		
Average sentiment level	Module of average sentiment level of tweet obtained from Sentiment WordNet		
Average sentiment score	Average sentiment level of tweet obtained from Sentiment WordNet		
Number of misspellings	Number of misspellings fixed by Microsoft Word spell checker		
Number Of Mistakes	Number of words that contains mistake but cannot be fixed by Microsoft Word spell checker		
Number of rejected tweets	Number of tweets where 70% of words either not in English or cannot be fixed by Microsoft Word spell checker		
Number of terms average	Average number of terms per / tweet		
Number of Foursquare check- ins	Number of Foursquare check-ins performed by user		
Number of Instagram medias	Number of Instagram medias posted by user		
Number of Foursquare tips	Number of Foursquare Tips that user post in a venue		
Average time between check- ins min	Average time between two sequential check-ins - represents Foursquare user activity frequency		

- Location features
 - Location semantics
 - Location topics

We map all Foursquare check – ins to Foursquare categories from category hierarchy.

For case when user performed check-ins in two restaurants and airport but did not perform check-ins in other venues:

	Category ₁		Category _{restaurant}		Category _{airport}		Category _n
U ₁	0	0	2	0	1	0	0
	*	*	*	*	*	*	*
U _n	*	*	*	*	*	*	*

Venue semantics such as venue categories can be related to users demography. E.g. individuals who tent to visit night clubs are usually belong to 10 – 20 or 20 – 30 years old age groups.

- Image features
 - Image concept learning

Extracted image concepts may represents user interests and be related to one's demography. For example female user may take pictures of flowers, food, while male – of cars or buildings.

*The concept learning Tool was provided by Lab of Media Search LMS.

It was evaluated based on ILSVRC2012 competition dataset and performed with average accuracy @10 - 0.637

Ensemble learning

Ensemble learning

Ensemble learning details

- According to our evaluation, the bias of estimated ages does not exceed ±2.28 years. It is thus reasonable to use the estimated age for age group prediction task.
- We have adopted SMOTE* oversampling to obtain balanced age-group labeling
- By performing 10-fold cross validation, we determine the optimal number of constructed random trees for each classifier with iteration step equal to 5 as 45, 25, 35, 40, 105 random trees for Random Forest Classifiers learned based on location, LIWC, heuristic, LDA 50, and image concept features respectively.
- We jointly learn the l_i model "strength" coefficient by performing "Hill Climbing" optimization* * with step 0.05. The randomized "Hill Climbing" approach is able to obtain local optimum for non-convex problems and, thus, can produce resolvable ensemble weighting.

*N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. Smote: synthetic minority over-sampling technique. Journal of artificial intelligence research, 2002.

**An iterative algorithm that starts with an arbitrary solution to a problem, then attempts to find a better solution by incrementally changing a single element of the solution. If the change produces a better solution, an incremental change is made to the new solution, repeating until no further improvements can be found.

Experimental results (Singapore)

Method	Gender	Age		
State-of-the-arts techniques				
SVM Location Cat. (Foursquare)	0.581	0.251		
SVM LWIC Text(Twitter)	0.590	0.254		
SVM Heuristic Text(Twitter)	0.589	0.290		
SVM LDA 50 Text(Twitter)	0.595	0.260		
SVM Image Concepts(Instagram)	0.581	0.254		
NB Location Cat. (Foursquare)	0.575	0.185		
NB LWIC Text(Twitter)	0.640	0.392		
NB Heuristic Text(Twitter)	0.599	0.394		
NB LDA 50 Text(Twitter)	0.653	0.343		
NB Image Concepts(Instagram)	0.631	0.233		
Single-Source				
RF Location Cat. (Foursquare)	0.649	0.306		
RF LWIC Text(Twitter)	0.716	0.407		
RF Heuristic Text(Twitter)	0.685	0.463		
RF LDA 50 Text(Twitter)	0.788	0.357		
RF Image Concepts(Instagram)	0.784	0.366		
Multi-Source combinations				
RF LDA + LIWC(Late Fusion)	0.784	0.426		
RF LDA + Heuristic(Late Fusion)	0.815	0.480		
RF Heuristic + LIWC (Late Fusion)	0.730	0.421		
RF All Text (Late Fusion)	0.815	0.425		
RF Media + Location (Late Fusion)	0.802	0.352		
RF Text + Media (Late Fusion)	0.824	0.483		
RF Text + Location (Late Fusion)	0.743	0.401		
All sources together				
RF Early fusion for all features	0.707	0.370		
RF Multi-source (Late Fusion)	0.878	0.509		

AGE GROUPS: < 20 YEARS OLD, 20 – 30 YEARS OLD, 30 – 40 YEARS OLD, > 40 YEARS OLD

Demographic mobility

User profile: Mobility + Demography

Geographical user mobility: users movement (city level)

Geographical user mobility: users movement (city level)

- Singapore population is concentrated in several regions, which represent peoples' housing (Regions 2 and 3) and working (Region 3) areas.
- There are some regions where male (Blue markers) user check-in density is much higher than female (Pink markers).

Geographical user mobility: users movement (region level)

Geographical user mobility: users movement (region level)

- Both female and male users often perform trips to nearby cities for shopping and leisure purposes (Regions 1, 2, 4, 5).
- Regions 2 and 3 are popular among female users, since 2 is "Malacca resorts", while 3 – National park. Both regions are famous by it's family time spending facilities.

Geographical user mobility: users movement (city level)

Geographical user mobility: users movement (city level)

- Teenagers and children (Brown markers) mostly perform check-ins in housing city areas and around schools (Regions 1,2,3,5).
- Students (Green markers) and working professionals (Blue and Red markers) are concentrated in city center (Region 4).

Geographical user mobility: users movement (region level)

Geographical user mobility: users movement (region level)

- Young users (brown circles) are rarely travel to nearby cities due to their age (Region 3)
- Adults (green circles) often make such trips (Regions 1 and 2). These users may be students or young professionals who visit their families during weekends.

Dataset Statistics: Content

Geographical user mobility: venue semantics profiling

 We extract location topics based on venue categories to model user mobility semantics

LDA word distribution over 6 topics for collected Foursquare check-ins.

Location topics may serve as an user interest clusters for distinguishing user demography attributes such as age or gender.

Every venue category

Table 2:	Category	distribution	among	LDA	topics
----------	----------	--------------	-------	-----	--------

ID	Categories	LDA Topics
T1	Malay Res-t, Mall, University, Indian	Food Lovers
	Res-t, Aisian Res-t	
T_2	Cafe, Airport, Hotel, Coffee Shop,	Travelers
	Chinese Res-t	(Business)
T3	Nightclub, Mall, Food Court, Trade	Party Goers
	School, Res-t, Coffee Shop	
T4	Home, Office, Build., Neighbor-d,	Family Guys
	Gov. Build., Factory	(Youth)
T5	University (Collage), Gym, Airport,	Students
	Hotel, Fitness Club	
T6	Train St., Apartment, Mall, High	Teenagers
	School, Bus St.	(Youth)

Geographical user mobility: venue semantics profiling

Geographical user mobility: venue semantics profiling

- Male users more often do shopping than male, while female users often show-up in job-related venues.
- > 30 years old users often show-up in dining-related places, while < 20 – often visit education-related venues.

Future work

Future work: Extended User Profiling

- Extended Demographic Profiling:
 - Occupation detection;
 - Personality detection;
 - Social status detection.
- Extended Mobility Profiling :
 - User communities detection and profiling (In terms of demographics, movement patterns, multi-source interests);
 - Cross-region mobility profiling (comparison of users' mobility across different regions and cultures).

Future work: Sensor Data Incorporation & Wellness Research

> Wellness lifestyle recommendation via:

- Chronic diseases tendency prediction
- Cross-source causality relationships analysis (just like Ramesh Jain proposed*)

Future work: How the framework may look like

Other task based on NUS-MSS

- **1.** Demographic profile learning
- 2. Multi-source data fusion
- **3.** Individual and group mobility analysis
- 4. Cross-source user identification
- 5. Cross-region user community detection
- 6. Cross-source causality relationships extraction
- 7. Users' privacy-related and cross-disciplinary research

Conclusions

- 1. We constructed and released a large multisource multi-modal cross-region "NUS-MSS" dataset;
- 2. We conducted first-order and higher-order learning for user mobility and demographic profiling;
- 3. New multi-modal features were proposed for a demographic profile learning.
- 4. Based on our experimental results, we can conclude that multi-source data mutually complements each other and their appropriate fusion boosts the user profiling performance.

Thank you!

You could download NUS-MSS dataset from: <u>http://lms.comp.nus.edu.sg/research/</u> **NUS-MULTISOURCE.htm** OR http://nusmultisource.azurewebsites.net **Aleksandr Farseev** National University of Singapore e-mail: farseev@u.nus.edu

* Ground truth construction

* Multi-source user Id mapping

* Text preprocessing

- Retweet filter filters out all retweeted tweets since it does not bring any information about users demography i.e. posted by other user;
- > Hash tags filter filters out all hash tags from user tweets;
- Slang transformation filter transforms all slang words to synonyms from dictionary;
- User mentions and place mentions filter filters out all user and place mentions;
- Repeated chars transformation filter filters out all repeated characters from tweets.