i-ASC 2014 Workshop in conjunction with ECIR 2014, Amsterdam, the Netherlands, 13th April 2014

# Mining digital footprints for smart tourism

Raffaele Perego, ISTI-CNR, Italy

I. Brilhante, C. Lucchese,  J.F. de Macedo, F.M. Nardini, C. Renso, C. Muntean, F. Silvestri, U. Vespier

**UCG**

**Knowledge**
- PoIs, Poi popularity, trends, patterns
- Signals of collective/personal behaviors
- How the above change with demography, seasonality, events

**Applications**
- Travel and Tourism Market
- Cultural Heritage fruition analysis
- Tourism and City governance

# Twitter

- ~1% of tweets are geotagged
- ~25% of geotagged tweets are LBS checkins
- Easy to crawl via the public API
  - 1% but can access more if you ask for a specific topic or place

# Facebook

- Huge repository of user personal information
- Friendship constraints on posts visibility
  - Large crawl not possible
- FB Apps can collect data
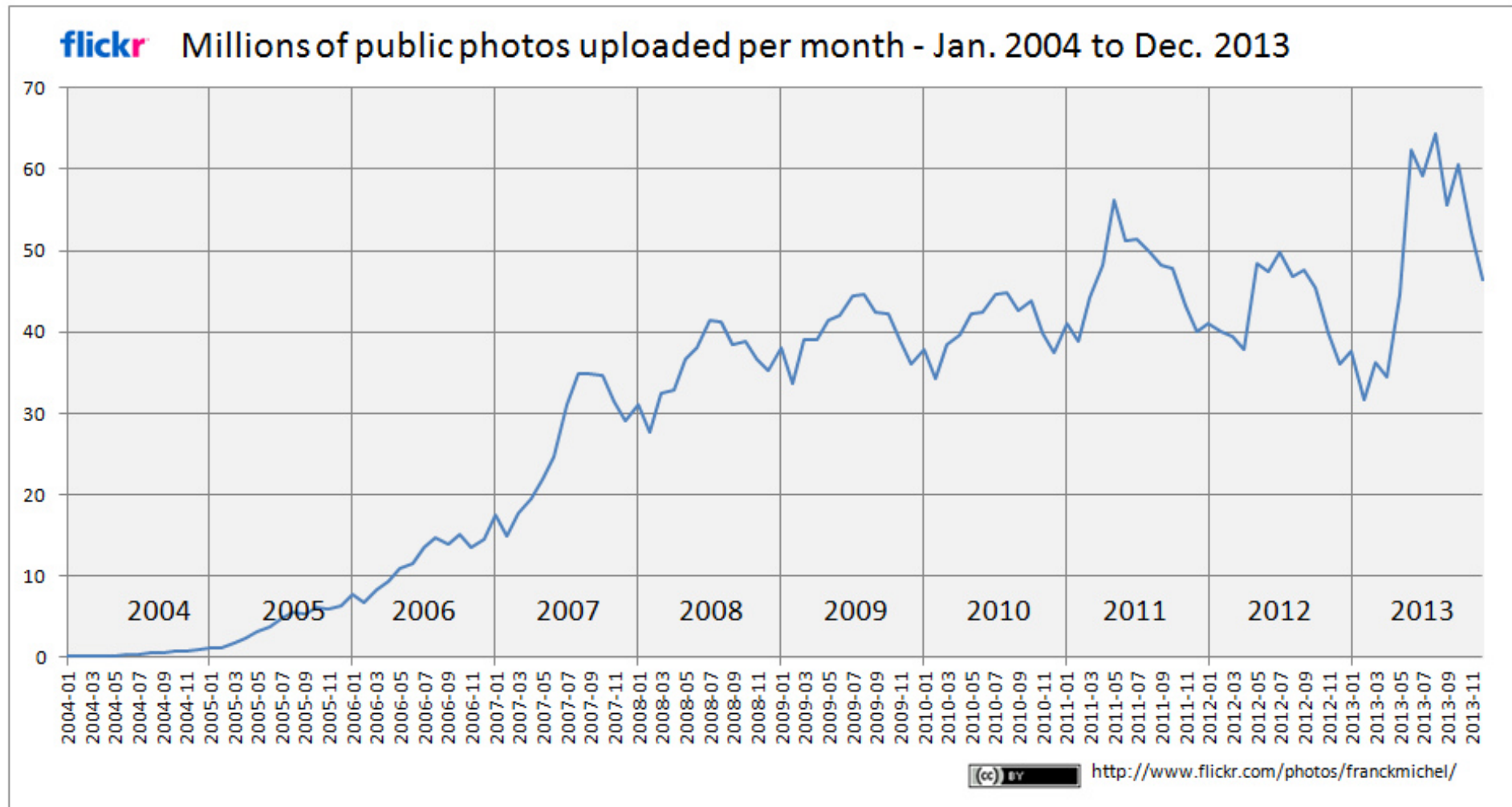
# Location-Based Services

*I'm at Stazione Roma Termini - @fsnews_it (Roma, RM) w/ 12 others http://4sq.com/1l4mk81*

- Users can check-in at a pre-defined place, or create a new place
  - Data biased toward sponsored listings
- Not available for public use:
  - Crawl Twitter as people can choose to publish their Foursquare check-ins on Twitter.

# Flickr

- Vast amount of rich data
  - 586 M public Photos uploaded in 2013
  - (Geo-)Tags, Titles, likes, Descriptions, Comments, Social profiles
- Easy to crawl
- Existing large public crawls:
  - CoPhIR: http://cophir.isti.cnr.it/
- Bulk uploading very common

# How many photos are uploaded to Flickr every day, month, year?



flickr  Millions of public photos uploaded per month - Jan. 2004 to Dec. 2013

edinburghcastle edinburgh

exchangesquare manchester

trafalgarsquare london

damsquare amsterdam

copenhagen copenhagen

cathedral köln

brandenburggate berlin

oconnellstreet dublin

bathabbey bristol

eiffel paris

praçadocomércio lisbon

plazamayor madrid

sagradafamilia barcelona

casino monaco

galleria milano

coliseum rome

pontevecchio firenze

sanmarco venice

rathaus münchen

europe praha

(credits to David Crandall et al., Cornell University)

# Our Challenges


pontevecchio trip firenze


palazzo vecchio canon florence

Given a large and noisy collection of photo albums taken in a given city:

1. Clean and organize the collection in semantically coherent clusters

2. Associate relevant PoIs with these clusters

3. Devise routes of tourists through these PoIs and characterize as precisely as possible their behaviors

4. Exploit such knowledge to provide personalized recommendations

- Where shall we go today?: planning touristic tours with Tripbuilder. CIKM 2013.

- TripBuilder: A Tool for Recommending Sightseeing Tours. Demo. ECIR 2014.

- LearNext: learning to predict tourists movements. CIKM 2013

Locals and tourists
(credits to Erik Fisher)

Amsterdam

**flickr**

**Match Photos to PoIs**

**Colosseum**
3 photos
01/07/2013 9:00 -12:00

**Ruins**
2 photos
01/07/2013 13:30 -15:00

**Trevi Fountain**
2 photos
01/07/2013 15:42 - 16:00

**Devise patterns of tourists behavior**

# Solution 1:
# exploit visual content

# Flickr Geo-Tags in Florence



How do we spatially group the photos?

# Geo-clustering with DBSCAN

# Visual Clustering



geocluster's images

Represents

01001...0
00111...1

.
.
.

11010...1

k-dimensional visual vectors

local descriptors

$x_{11}$  $x_{12}$  $x_{1m}$

$x_{n1}$  $x_{n2}$  $x_{nm}$

...

$c_{11}$  $c_{12}$  $c_{1k}$

$c_{n1}$  $c_{n2}$  $c_{nk}$

k visual words (centroids)

K-Means Clustering

Goal: to reduce the cost of computing similarity

H. Bay, T. Tuytelaars, and L. Van Gool. Surf: Speeded up robust features. Computer Vision–ECCV 2006
G Csurka, C. Dance, L Fan, J Willamowski, and C. Bray. Visual categorization with bags of keypoints. ECCV, 2004.

# Labeling with tags



florence firenze dmclx2 torres lumix
f l o r e n z   t o w e r s        i t a l i e n
palazzovecchio tuscany leica tours
architecture digital italy landscapes
palazzodellasignoria torri

Two key ideas:

- ## Using the *spatial relevance* of tags
  - Measure: ratio between the tag area and the overall geographical area analyzed

- ## Using the *social relevance* of tags
  - Measure: number of different users using a given tag

$$\mathrm{GEORELEVANCE}(tag_k) < t_{geo} \wedge \mathrm{SOCIALRELEVANCE}(tag_k) > t_{social}$$

# Simple Demo

http://hpc.isti.cnr.it:8000

Flickr photos collection of Florence

| | |
|---|---|
| Images crawled | 53563 |
| Geo-clusters | 112 |
| Geo-clustered images | 37187 |
| Visual Clusters | 743 |
| Visual clustered images | 4235 |

# Solution 2:
## using geo-tags and Wikipedia

# Trajectories from Flickr & Wikipedia

flickr

WIKIPEDIA
The Free Encyclopedia

**Colosseum**
3 photos
01/07/2013 9:00 -12:00
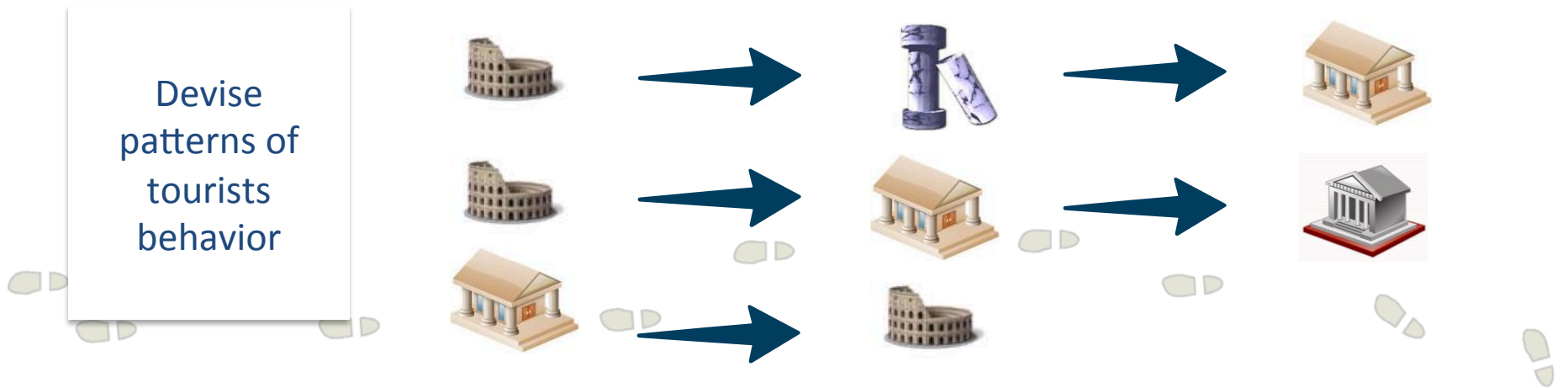
**Ruins**
2 photos
01/07/2013 13:30 -15:00

**Trevi Fountain**
2 photos
01/07/2013 15:42 - 16:00

Devise
patterns of
tourists
behavior

# Datasets

| City | PoIs | Users | Photos | Trajectories | Traj. per PoI (avg.) |
|---|---|---|---|---|---|
| Pisa | 112 | 1,825 | 18,170 | 3,430 | 7.20 |
| Florence | 891 | 7,049 | 102,888 | 16,522 | 5.39 |
| Rome | 490 | 13,772 | 234,616 | 35,522 | 20.51 |

# Wikipedia Categories

- Categories for Colosseum:
  - 1st-century architecture
  - Amphitheatres in Rome
  - Ancient Roman architecture
  - Building projects of the Flavian dynasty
  - Roman archaeology
  - Ruins in Italy

# User and Trajectory profiling

*Categories: Amphitheatres in Rome, Ancient Roman architecture, Roman archaeology, Ruins in Italy*



- A profile for a user can be built from the categories of the PoIs visited
  - We can measure the "interest" of a user for a (set of) PoI, e.g.:

$$\Gamma(p, u) = \alpha \cdot sim(\vec{v_p}, \vec{v_u}) + (1 - \alpha) \cdot pop(p)$$

- Each trajectory can be labeled with the set of categories of the constituent PoIs

# Planning Sightseeing Tours with TripBuilder



What should I visit in San Francisco?

Given:
- Time: 2 days;
- My preferences

**Golden Gate Bridge**

**Golden Gate Park** — 4 h

**California Academy of Sciences**

**de Young Museum** — 4 h

**San Francisco Museum of Modern Art** → **Aquarium of the Bay** → **Alcatraz** — 8 h

How many of these "trajectories" visit such places?

How do users compose tours?

# The TripCover Problem

- **Given:**
  - A set of popular trajectories crossing a set of PoIs and their time cost
  - The relevance of the trajectories w.r.t. the category set
  - The **Time Budget** and **Preferences** of a user
  - A measure of **PoI-User interest**

- **Find:**
  - the subset of trajectories that maximizes user interest and fits in the time budget



**TripCover** is an instance of the **Generalized Maximum Coverage** (GMC) problem. NP-Hard with a (e/(e-1))-approximation algorithm.

# TrajSP: joining the trajectories

- TripCover solution is a set of trajectories fitting user interest and time budget
    - Local search heuristics based on 2-opt and 3-opt moves for connecting the solution in a single sightseeing tour

# Next PoI prediction with LearNext



Learn the "next" PoI

Learning to rank
- Ranking SVM
- GBRT

Feature engineering

# Features

| Popularity | Frequent seq. of PoIs | User preference | Distance and time | Poi Characteristics | Session Characteristics |
|---|---|---|---|---|---|

## Session features:

- actualTransferTime
- actualVisitTime
- distLat_Avg/Max/Min/Tot
- distLen_Avg/Max/Min/Tot
- euclideanDist_Avg/Max/Min/Tot
- sessTime
- phPoISess_Avg/Max/Min/Tot
- categsPerSess
- uniqueCategsPerSess
- sessLen
- userSessLen_Avg/Max/Min/Tot
- userSessRatio

## Candidate PoI features:

- distFromFirstPoI_Eucl/Lat/Len
- distFromLastPoI_Eucl/Lat/Len
- visitTimePoI_User
- visitTime_Avg/Max/Min/StDev/Tot
- freqBigrams
- freqTrigrams
- start/stop/middleProbab
- cat1, cat2, ..., cat10
- entropy
- numCategs
- numPhotos_Avg/Max/Min/Tot
- noOfVisists
- photosPoIUserId_Avg/Total
- ratioPhotosPoI
- ratioSessWithPoI
- ratioUsersVisitingPoI
- photosPerUser
- ratioPoIInUserPhotos

# Learning to Rank approach

- Building a model that ranks highest the PoI most likely to be visited as next by the tourist.
  - A trail is represented by a 68-dimension feature space.
  - Each example is represented by the feature vector and its label indicates the PoI's degree of relevance to the user {0,1}.

- The learning algorithm is trained to predict the relevance from the feature vector. Approaches:
  - Ranking SVM
  - GBRT

# Experiments

- ## Baselines:
  - Probability: suggesting the most probable next PoI from the current one
  - WhereNext[1]: a trajectory pattern mining approach
  - Random Walk[2]: a graph based approach "Itinerary Graph" exploiting RWR

- ## Training/test methodology
  - Training: 1 positive, 3 negative
  - Test: all unseen PoIs for the trail

- ## Metrics:
  - Success@k
  - MRR@k
  - MRR

| | Session | Candidate PoI | Relevance |
|---|---|---|---|
| | | ✓ | Positive Example |
| Feature vector : | $\{fs_1, fs_2, ..., fs_n\}$ | $\{fp_1, fp_2, ..., fp_m\}$ | 1 |
| | | ✗ | Negative Example |
| Feature vector : | $\{fs_1, fs_2, ..., fs_n\}$ | $\{fp_1, fp_2, ..., fp_m\}$ | 0 |

1. A. Monreale, F. Pinelli, R. Trasarti, and F. Giannotti. "WhereNext: a location predictor on trajectory pattern mining". In Proc. SIGKDD. ACM, 2009.
2. C. Lucchese, R. Perego, F. Silvestri, H. Vahabi, and R. Venturini. "How random walks can help tourism". In Proc. ECIR. LNCS, 2012.

# Results

| City | Predictor | Success (MRR) | | | | | MRR |
|------|-----------|------|------|------|------|------|-----|
| | | @1 | @2 | @3 | @5 | @10 | |
| Pisa | PROB | 16.08% | - | - | - | - | - |
| | WhereNext [11] | 12.56% | - | - | - | - | - |
| | Random Walk [10] | 15.07% (0.15) | 20.60% (0.17) | 25.12% (0.19) | 31.65% (0.20) | 46.73% (0.22) | - |
| | Ranking SVM | 32.66% (0.32) | 49.74% (0.41) | 55.77% (0.43) | 65.82% (0.45) | 73.36% (0.46) | 0.47 |
| | GBRT | 40.70% (0.40) | 55.27% (0.47) | 63.81% (0.50) | 75.87% (0.53) | 88.44% (0.55) | 0.56 |
| Florence | PROB | 4.59% | - | - | - | - | - |
| | WhereNext [11] | 2.90% | - | - | - | - | - |
| | Random Walk [10] | 3.25% (0.03) | 6.09% (0.04) | 8.77% (0.05) | 11.69% (0.06) | 20.13% (0.07) | - |
| | Ranking SVM | 33.91% (0.33) | 41.01% (0.37) | 44.27% (0.38) | 48.20% (0.39) | 53.29% (0.40) | 0.41 |
| | GBRT | 37.76% (0.37) | 46.78% (0.42) | 53.04% (0.44) | 59.31% (0.45) | 69.34% (0.47) | 0.48 |
| Rome | PROB | 12.93% | - | - | - | - | - |
| | WhereNext [11] | 6.37% | - | - | - | - | - |
| | Random Walk [10] | 8.43% (0.08) | 13.76% (0.11) | 19.22% (0.12) | 26.38% (0.14) | 38.12% (0.16) | - |
| | Ranking SVM | 21.88% (0.21) | 30.24% (0.26) | 36.37% (0.28) | 46.95% (0.30) | 59.49% (0.32) | 0.33 |
| | GBRT | 30.95% (0.30) | 40.07% (0.34) | 47.15% (0.38) | 56.34% (0.40) | 67.68% (0.41) | 0.42 |

Print-outs of all the images uploaded to Flickr in a day
(installation by Erik Kessels, Amsterdam)



THANKS!