



Aalto University
School of Science
and Technology

Exploring San Francisco's geography using geo-tagged photos

Aristides Gionis Michael Mathioudakis
Géraud Le Falher

December 10, 2013

Outline

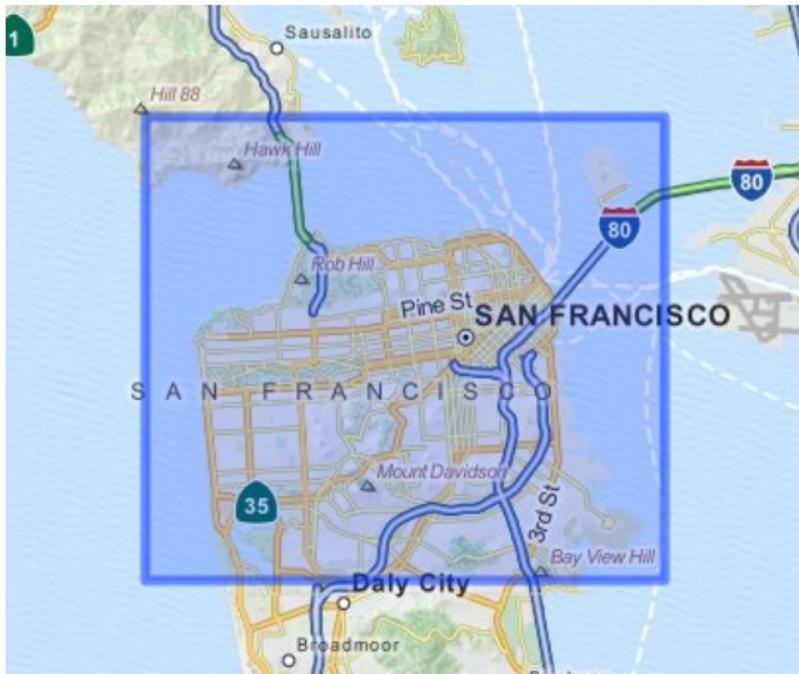
Introduction

Methods

Result

Future work

Dataset



San Francisco: 780K photos and 140K tags since January 2008

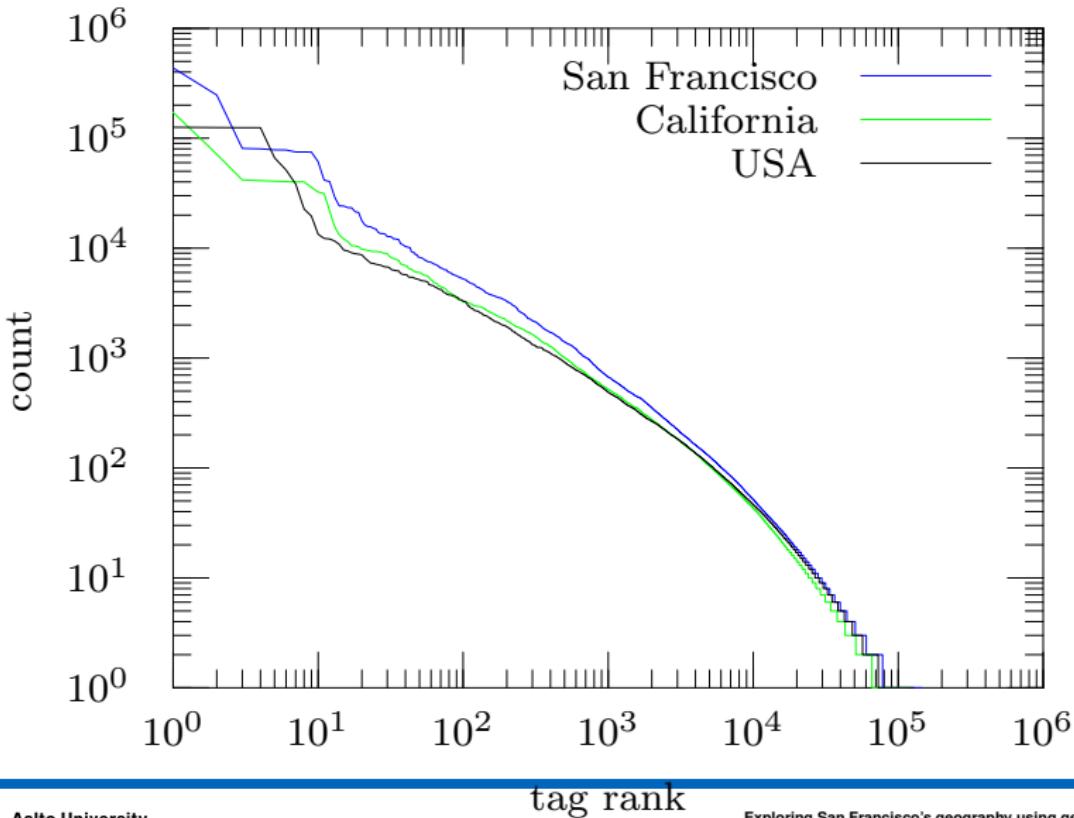
Photo metadata

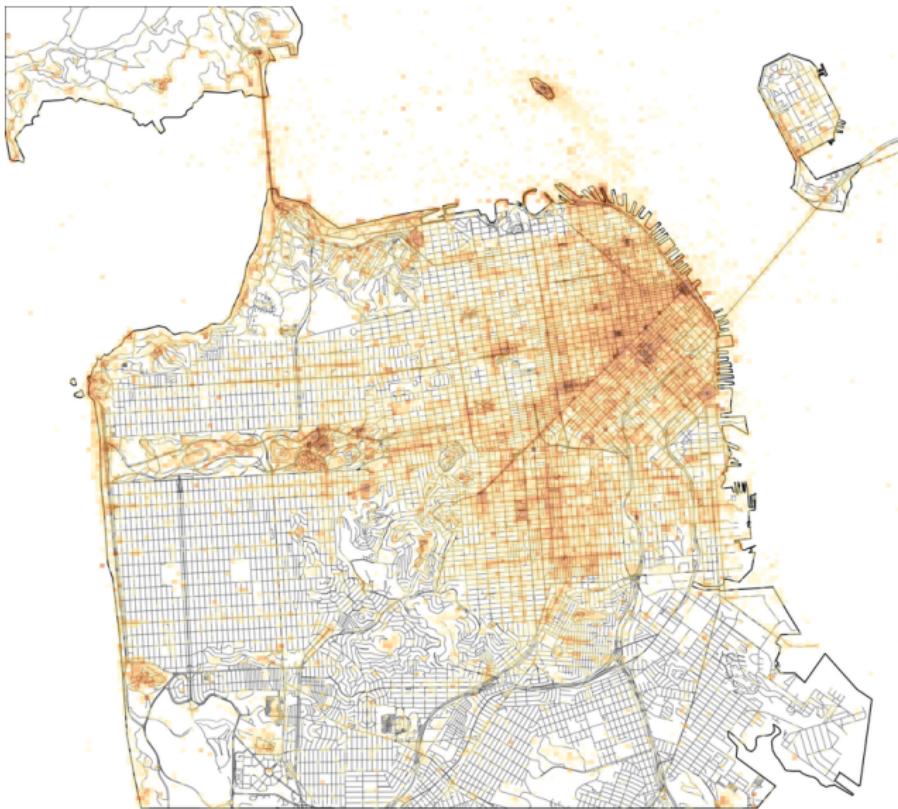
```
"loc"      : [-122.392501, 37.77515],  
"taken"    : "2008-03-24 14:55:40",  
"user_id"  : "37417902@N00",  
"tags"     : ["sanfrancisco", "california",  
             "bridge", "chinabasin"]
```

Problem

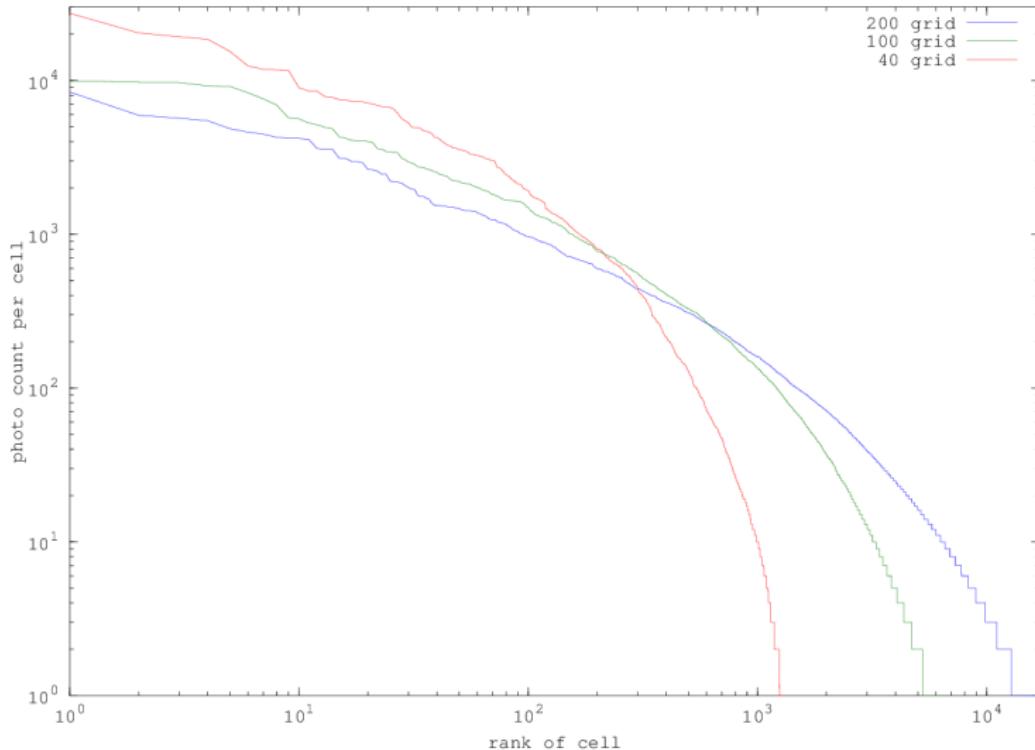
- ▶ The goal is to find association rules
 - ▶ tag: in which location it appears
 - ▶ locations: which tags describe it
 - ▶ city: which locations are interesting
- ▶ Applications
 - ▶ Flickr website
 - ▶ Tourism

Exploration of data





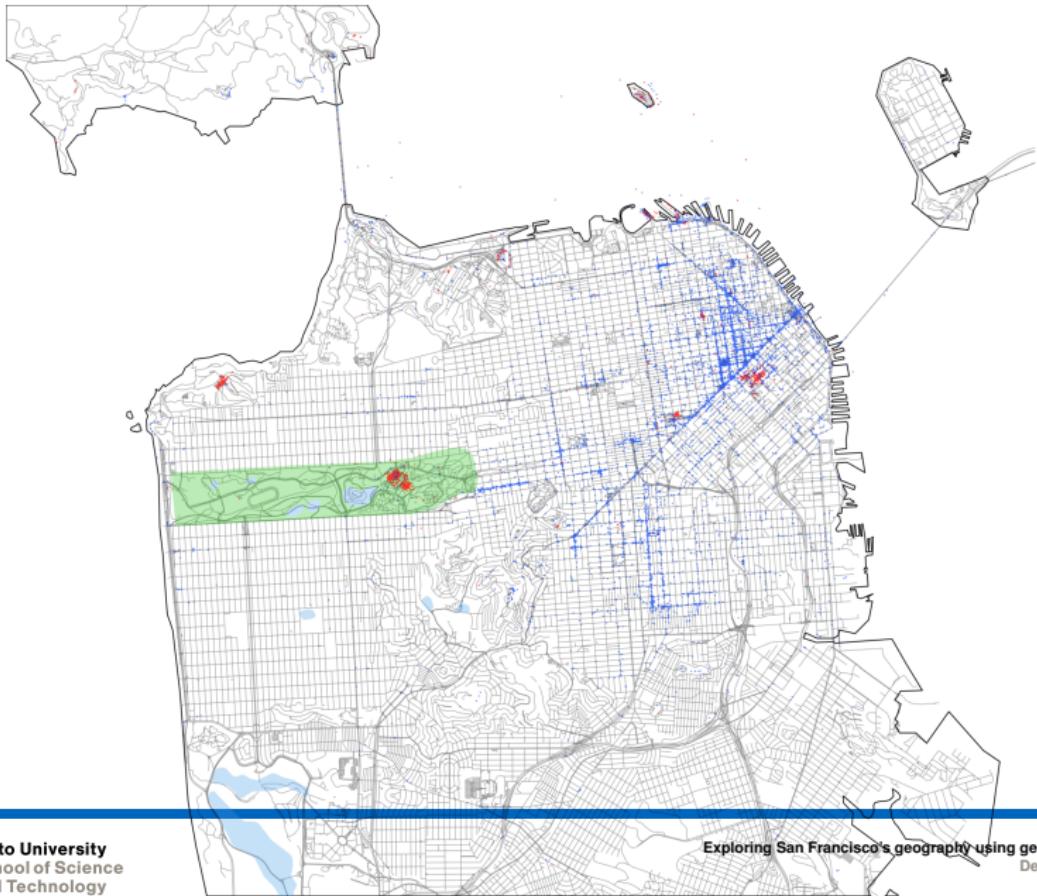
Photos distribution



Preliminary cleaning

- ▶ Avoid spurious photos by removing tags of the same user set almost in the same place at almost the same time.
- ▶ Keep only tags used enough times by enough users over enough time.

Statistics



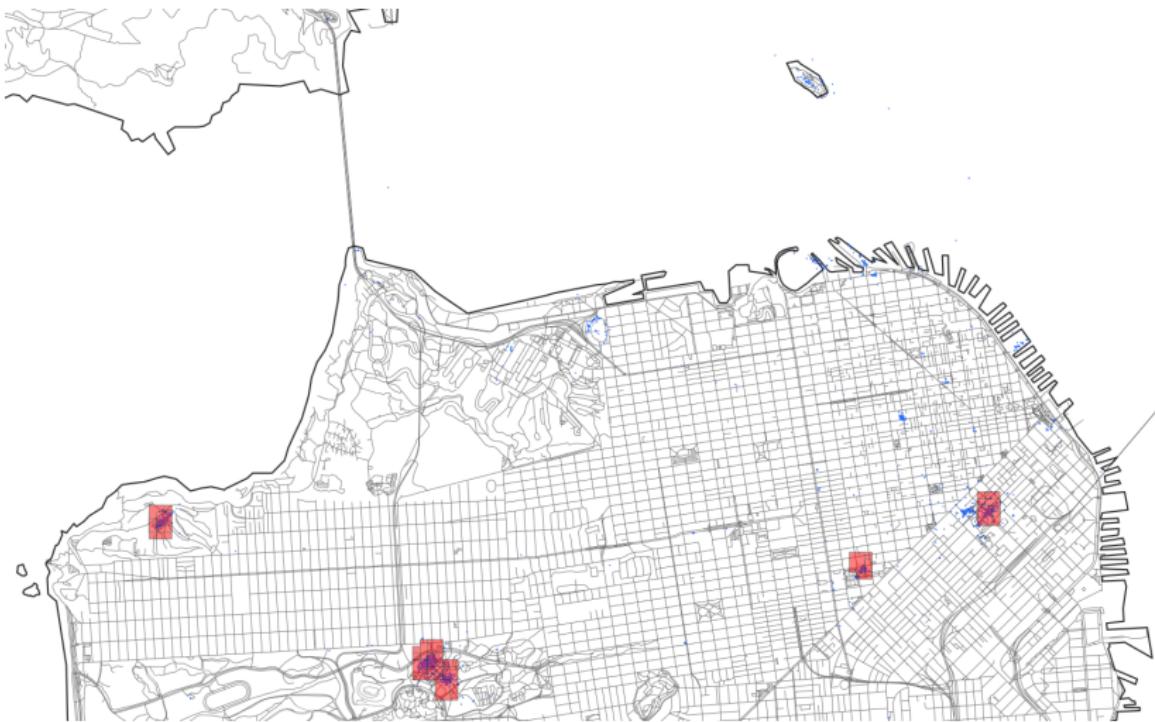
- ▶ First, Entropy and Kullback Leibler divergence with background photos:

$$H(\text{museum}) = 3.565 \quad H(\text{street}) = 6.875$$

Lowest entropy	Highest divergence
californiaacademyofsciences	zoo
conservatoryofflowers	treasureisland
sfmoma	conservatoryofflowers
deyoung	oceanbeach
deyoungmuseum	ucsfschoolofdentistry
ucsfschoolofdentistry	ucsf
attpark	japantown
cityhall	sfmoma

- ▶ But it provide only one value for each tag, i.e. it loses spatial information
- ▶ Enter Kulldorff Spatial Scan Statistic: for a region R , compute the KL divergence considering only R and its complementary: $t(R) \log \frac{t(R)}{b(R)} + (1 - t(R)) \log \frac{1-t(R)}{1-t(R)}$
- ▶ Computed for all supported tags over all reasonably sized regions

Where are the museums?



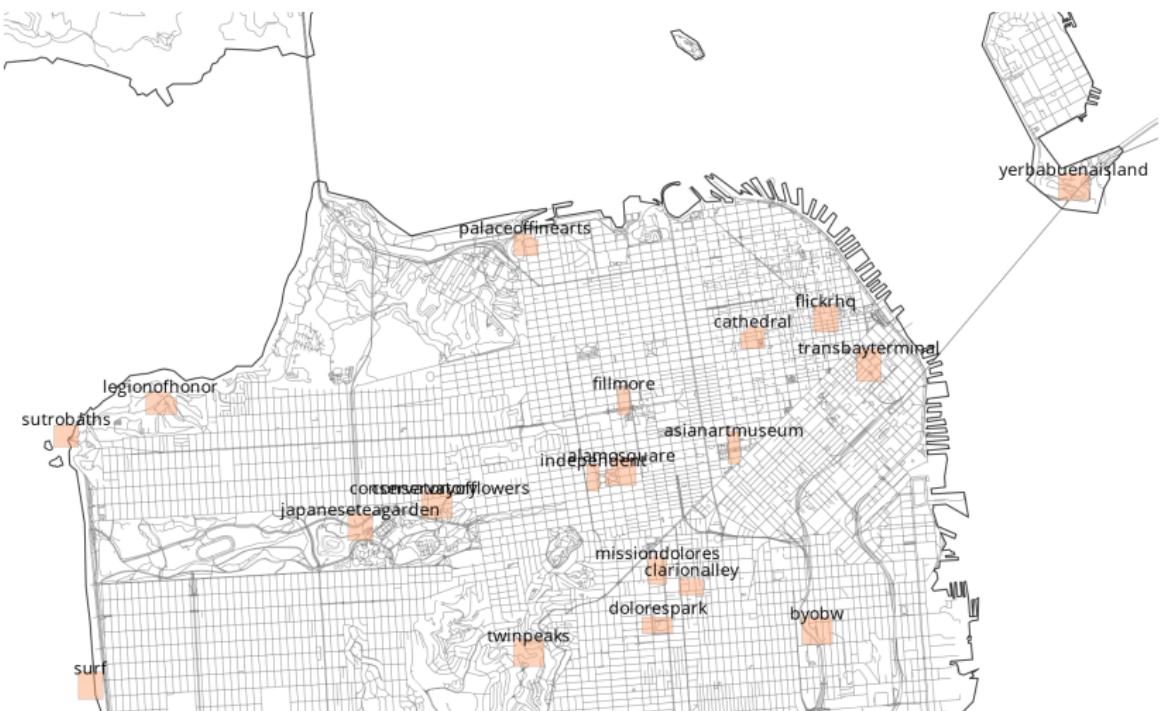
What is this place?







Where to go in San Francisco?



Future work

- ▶ Evaluate result using ground truth: tourist guide, manual annotation
- ▶ Use less parameters and more information from the data
- ▶ Extend to different scales: city, region, country, world
- ▶ Consider more type of entities: time, users



Thank you

Questions?