



Figure 1: Illustration of the RRS algorithm. The background color represents venue densities as shown on the colorbar. (a) Two starting points $l1_{sw}$ and $l1_{ne}$ are assigned with an initial venue density d_0 . Hence, the predicted density for $l1_1$ is $d'_1 = d_0$, and the side length s_1 satisfies $s_1^2 = \bar{b}/d'_1$. By exhaustively searching this area, we learn its true venue density d_1 . (b) $l1_2$ is closer to $l1_1$ and $l1_{ne}$. Its density is predicted as the weighted average of d_1 and d_0 , with weights as the inverse of its distances to $l1_1$ and $l1_{ne}$. (c) $l1_5$ is chosen between two nearby sampled boxes, and its box is cut to be non-overlapping with them. (d) Over steps, a sequence of sample boxes are formed with more and more accurate densities predicted. Locations falling into already sampled boxes trigger the corresponding boxes being re-sampled.

However, it is very costly to perform such an exhaustive search on a large and dense geographic region, for example, New York City, where the dimension of small searching boxes needs to be within a few meters to meet the criterion. Hence, this method is not suitable for timely learning and estimating the statistics of LBSNs. This motivates us to design efficient sampling methods to collect a representative venue set from the objective region and design unbiased estimators for estimating the statistics such as the total number of venues and venue label distributions.

3. RANDOM REGION SAMPLING

The basic idea behind the random region sampling (RRS) is that given a sampling budget B , i.e., the number of queries allowed, at every step, a location, $l1$, is chosen uniformly at random from the objective region G , and the size of the sample box around $l1$ is determined by two criteria as follows.

- **Box size selection using venue density prediction:** The venue density d' around $l1$ is predicted as the weighted average of the venue densities of its closest sampled boxes. The side length s of a new box centered at $l1$ is computed as $\sqrt{\bar{b}/d'}$, to keep the expected number of venues in this new box close to the API return limit \bar{b} .
- **Non-overlapping boxes:** Check whether the new box obtained above collides with any already sampled box, and cut the new box if necessary to keep it containing $l1$ and non-overlapping with those sampled boxes.

Using the above two criteria, a non-overlapping new box G_{l1} is determined based on the best knowledge to have an expected total number of venues close to the API return limit. Then, an exhaustive search is performed on this box, which cumulates one true venue density in G , thus improves the following venue density predictions. The area it covers is considered as a sampled box. Later if a random location is chosen that falls into this box, it is considered that this box is sampled again, and no actual API queries are needed. Note that keeping the sampled boxes to be non-overlapping ensures each box having an invariant probability to be sampled again once the box has been established in one sampling process. This probability is proportional to the size of its area, and it is an important quantity in designing unbiased estimators for the objective region. Until running out of the API budget B , m samples, X_1, \dots, X_m are drawn from n non-overlapped boxes G_1, \dots, G_n , with $m \geq n$. If the budget runs out while exhaustively searching the last box, that box will be ignored. Fig 1 takes New York City as an example to illustrate how RRS works.

Estimators. Theorem 1 below presents how to use the venues sampled with RRS to estimate the total number of venues N of objective region G , by the estimator \hat{N} , where we omit the proof and the evaluation results due to the space limit.

THEOREM 1. *Using RRS with budget B , we obtain m sampled boxes X_1, \dots, X_m . Let $a(X_t)$ and $f(X_t)$ be the size and number*

of venues of X_t , $1 \leq t \leq m$. Then, \hat{N} in eq.(1) is an (asymptotically) unbiased estimator of N .

$$\hat{N} = \frac{a(G)}{m} \sum_{t=1}^m \frac{f(X_t)}{a(X_t)}. \quad (1)$$

4. EXPLORING VENUE POPULARITY

Now, we investigate how various factors affect the venue popularity in Foursquare, by analyzing a unique dataset with 2,398,931 venues collected from 14 regions during 05/01/12–06/30/12, including New York City, Paris, Seoul, covering a wide range of geographic areas. Our main results are summarized as follows:

- **Venue profile.** Venues with more complete profile information are more likely to be popular, and the two most influential individual attributes are “contact” and “cross street”.
- **Venue category.** By performing comprehensive categorical analysis, we observe that venues in the Food category attract the most (43%) public comments (tips) by users, and the Travel & Transport category is the most popular category with the highest per venue check-ins. The residence, office, and school have the highest *user stickiness*, i.e., the average number of repeated visits of users to each venue.
- **Venue age.** The most popular venues were usually created at the early phase of Foursquare.

5. CONCLUSION

Taking Foursquare as an example, we study how to sample and estimate characteristics of location based social networks. Moreover, the findings of our venue popularity analysis may help advertisers to select promising candidate venues for effective advertisement placement, and venue owners to improve venues’ attraction to customers.

6. REFERENCES

- [1] Foursquare. <https://foursquare.com/>.
- [2] Foursquare API. <https://developer.foursquare.com/>.
- [3] Foursquare on wiki. <http://en.wikipedia.org/wiki/Foursquare>.
- [4] Jie Bao, Yu Zheng, and Mohamed F. Mokbel. Location-based and preference-aware recommendation using sparse geo-social networking data. In *ACM GIS*, 2012.
- [5] Abdelaziz Mohaisen, Pengkui Luo, Yanhua Li, Yongdae Kim, and Zhi-Li Zhang. Measuring bias in the mixing time of social graphs due to graph sampling. In *Milcom*, 2012.
- [6] Anastasios Noulas, Salvatore Scellato, Cecilia Mascolo, and Massimiliano Pontil. An empirical study of geographic user activity patterns in foursquare. *ICWSM'11*.
- [7] Jia Zhou, Yanhua Li, Vijay Kumar Adhikari, and Zhi-Li Zhang. Counting youtube videos via random prefix sampling. In *IMC*, 2011.