

Opportunistic Sampling for Joint Population Size and Density Estimation: Supplementary Notes

Farid M. Naini, Olivier Dousse, Patrick Thiran, and Martin Vetterli

Abstract

In this document we provide supplementary material for [1].

CONTENTS

I	Population Size Estimation Model	2
I-A	Poisson-Gamma Assumption of the Detection Patterns	2
I-B	Likelihood Function for the Population Size Estimation Model	3
I-C	Population Size Estimation by Using a Subset of the Information	3
II	Joint-Estimation Model for Population Size and Density	4
II-A	Likelihood Function for the Joint-estimation Model	4
III	Additional Measurements	6
III-A	Additional Measurements for the Paléo Dataset	6
III-B	Additional Measurements for the EPFL Dataset	6
III-B1	Users' Arrival and Departure Times	6
III-B2	Population Size Estimation using Subset of the Information	7
IV	Proof of the Theorems	8
IV-A	Background	8
IV-B	Proofs	8
	References	11

I. POPULATION SIZE ESTIMATION MODEL

A. Poisson-Gamma Assumption of the Detection Patterns

In Section 4.2.1 of [1] we demonstrate that the Poisson-Gamma model fits well to the contact patterns of individuals by agents. Here we show that the fit to the detection patterns is not good. We first redefine the Poisson-Gamma model based on the detection patterns.

We denote by κ_{ij} the number of times that individual i has been *detected* by agent j on the festival grounds. We denote by $\kappa_i = \sum_{j=1}^M \kappa_{ij}$ the total number of times that individual i has been detected on the festival grounds. Individual i is discovered (i.e., is among the S discovered individuals) if and only if $\kappa_i > 0$ (if he has been detected by at least one of the agents). The assumptions of our population size estimation model based on the detection patterns are as follows.

- **Poisson detections:** The number of times agent j detects individual i on the festival grounds, i.e., κ_{ij} , is Poisson distributed with mean equal to $\lambda_i t_{at_i, dt_i}^j$, where λ_i is called the *detection rate* of individual i .
- **Independence:** The random variable κ_{ij} is independent from $\kappa_{i'j'}$ for $i \neq i', j \neq j'$.

We assume that for individual i , λ_i is drawn from a Gamma distribution with unknown parameters α and β , independently from other individuals and from his arrival and departure times. Parameter λ_i represents how easily the individual *puts himself in a detectable position* on the festival grounds. We have

$$\kappa_{ij} \sim \text{Poisson}(\lambda_i \cdot t_{at_i, dt_i}^j).$$

We now check the fit of the above model to the detection patterns. Given the values of λ_i , at_i , and dt_i for individual i ,

$$\kappa_i = \sum_{j=1}^M \kappa_{ij} \sim \text{Poisson}(\lambda_i \sum_{j=1}^M t_{at_i, dt_i}^j).$$

That is κ_i is also Poisson distributed given λ_i , at_i , and dt_i . Now consider the values of $\kappa_1, \kappa_2, \dots, \kappa_S$ for the S discovered individuals; these values follow a *truncated* Poisson distribution, because the value of κ_i for individual i must be non-zero in order for the individual to be observed. The solid curve in Figure 1(a) shows the empirical distribution of the observed number of detections for all the S discovered individuals. As expected, number of detections of individuals is larger than their number of contacts, which is shown in Figure 4(a) of [1]. The dashed curve in Figure 1(a) shows the analytical distribution of the number of detections based on the truncated Poisson-Gamma

fit to the measurements. The corresponding Q-Q plot is shown in Figure 1(b), which shows that the fit of the Poisson-Gamma model to the measurements is not good, in particular they have very different tail behaviors. Pearson's chi-squared test for the equality of the two distributions gives a p-value of 8.61×10^{-9} , which indicates the rejection of the equality hypothesis of the two distributions (i.e., the empirical histogram of the detections and the corresponding fitted Poisson-Gamma distribution).

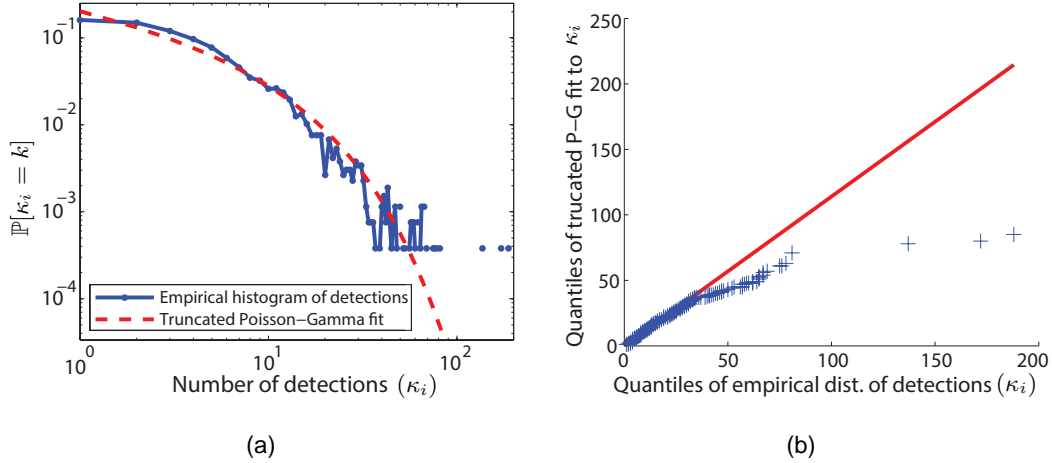


Fig. 1. The goodness of fit of truncated Poisson-Gamma distribution to the measurements. (a) and (b) show, respectively, the probability distribution function and the Q-Q plot with respect to the empirical distribution of detections.

B. Likelihood Function for the Population Size Estimation Model

In Section 4.3 of [1] we compute the likelihood function of our population size estimation model by using a distribution $f(at, dt)$ for individuals' arrival and departure times. In the case where in addition to $f(at, dt)$, individuals' exact arrival/departure times are also known, the likelihood function can be simplified as follows:

$$L(N, \alpha, \beta) = \binom{N}{S} \left(\mathbb{E}_{at, dt} \left[\left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at, dt}^j} \right)^\alpha \right] \right)^{N-S} \prod_{i=1}^S \left\{ \frac{\Gamma(\alpha + k_i) \beta^\alpha}{\Gamma(\alpha) \left(\beta + \sum_{j=1}^M t_{at_i, dt_i}^j \right)^{\alpha + k_i}} \right\}.$$

C. Population Size Estimation by Using a Subset of the Information

In Section 5.3 of [1] we apply our population size estimator to a subset of the Paléo measurements. In the first part, we consider the measurements obtained by a subset of size m of the agents. For each subset of size m , we consider all the possible combinations of the agents and

estimate the population size for each combination of m agents. The average percentage of the discovered individuals ($S/N \times 100$) and their 90% confidence intervals for all the combinations of size m are shown in Figure 2(a) for $m = 5, 6, \dots, 10$. In the second part, we consider the measurements obtained by all agents during an observation window of length w smaller than the festival duration. The observation window starting from the moment when the first agent arrives at the festival (17h09 for agent 5), until the moment when the last agent departs from the festival (4h01 for agent 1) is approximately 11 hours. We partition this interval into slots 10-minutes in length. For an observation window of length w , we consider all the consecutive 10-minute slots with total length w , and we estimate the population size based on the measurements obtained during these slots. The average percentage of the discovered individuals ($S/N \times 100$) and their 90% confidence intervals are shown in Figure 2(a), for $w = 4, 5, \dots, 11$ hours.

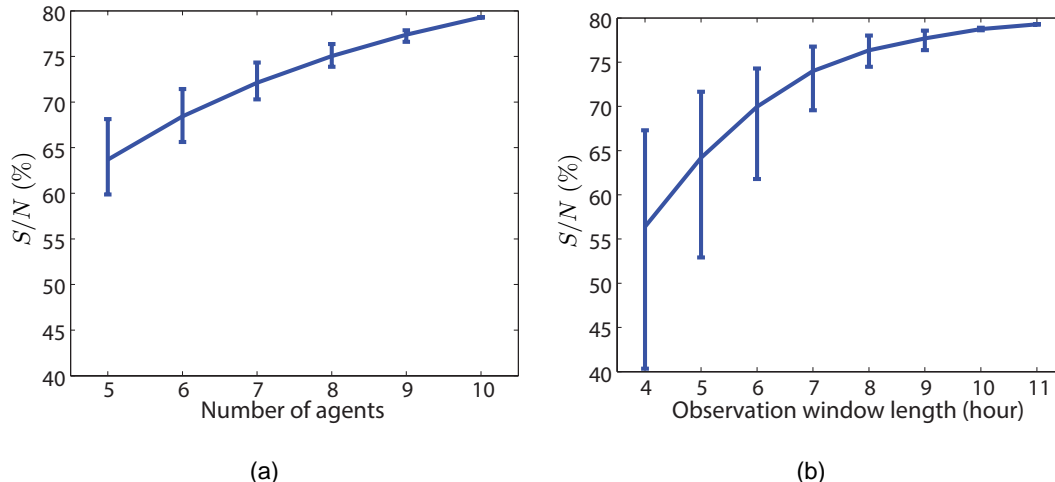


Fig. 2. Percentage of discovered individuals ($S/N \times 100$) as a function of (a) the number of agents and (b) the observation window length. The solid lines and the bars show the average and the 90% confidence interval, respectively.

II. JOINT-ESTIMATION MODEL FOR POPULATION SIZE AND DENSITY

A. Likelihood Function for the Joint-estimation Model of Population Size and Density

Here we derive the full likelihood function of the joint-estimation model of population size and density. The procedure for computing the likelihood function is similar to that of the population size estimation model given in Section 4.3 of [1].

Recall that we have the following property for the Gamma distribution:

$$\mathbb{E}_\lambda [e^{-\lambda x} \lambda^y] = \frac{\Gamma(\alpha + y) \beta^\alpha}{\Gamma(\alpha) (\beta + x)^{\alpha + y}}. \quad (1)$$

Given the location-dependent contact rates $\lambda^{(1)}, \dots, \lambda^{(K)}$, and the arrival/departure times at and dt for an individual, we have

$$p_{dsc}^{(\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(K)}, at, dt)} = 1 - \prod_{j=1}^M \prod_{l=1}^K e^{-\lambda^{(l)} t_{at, dt}^{j, (l)}} = 1 - \prod_{l=1}^K e^{-\lambda^{(l)} \sum_{j=1}^M t_{at, dt}^{j, (l)}}.$$

By marginalizing out the random contacts rates using Equation (1) we get

$$p_{dsc}^{(at, dt)} = 1 - \prod_{l=1}^K \left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)}} \right)^{\alpha^{(l)}},$$

and by marginalizing out the arrival/departure times of the individuals we get

$$p_{dsc} = 1 - \mathbb{E}_{(at, dt)} \left[\prod_{l=1}^K \left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)}} \right)^{\alpha^{(l)}} \right].$$

Now we compute the likelihood of the discovered individuals,

$$\begin{aligned} \mathbb{P}_i^{(\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(K)}, at, dt)} &= \prod_{j=1}^M \prod_{l=1}^K e^{-\lambda^{(l)} t_{at, dt}^{j, (l)}} \frac{(\lambda^{(l)} t_{at, dt}^{j, (l)})^{k_{ij}^{(l)}}}{k_{ij}^{(l)}!} \\ &= \prod_{l=1}^K e^{-\lambda^{(l)} \sum_{j=1}^M t_{at, dt}^{j, (l)}} (\lambda^{(l)})^{k_i^{(l)}} \prod_{j=1}^M \frac{(t_{at, dt}^{j, (l)})^{k_{ij}^{(l)}}}{k_{ij}^{(l)}!}. \end{aligned}$$

By marginalizing out the random contact rates using Equation (1) we get

$$\mathbb{P}_i^{(at, dt)} = \prod_{l=1}^K \frac{\beta^{\alpha^{(l)}} \Gamma(\alpha^{(l)} + k_i^{(l)})}{\Gamma(\alpha^{(l)}) (\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)})^{\alpha^{(l)} + k_i^{(l)}}} \prod_{j=1}^M \frac{(t_{at, dt}^{j, (l)})^{k_{ij}^{(l)}}}{k_{ij}^{(l)}!}.$$

By marginalizing out the arrival/departure times of the individuals we get:

$$\mathbb{P}_i = \mathbb{E}_{(at, dt)} \left[\prod_{l=1}^K \frac{\beta^{\alpha^{(l)}} \Gamma(\alpha^{(l)} + k_i^{(l)})}{\Gamma(\alpha^{(l)}) (\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)})^{\alpha^{(l)} + k_i^{(l)}}} \prod_{j=1}^M \frac{(t_{at, dt}^{j, (l)})^{k_{ij}^{(l)}}}{k_{ij}^{(l)}!} \right].$$

The total likelihood is:

$$\begin{aligned} L &= \binom{N}{N-S} \left(1 - \mathbb{E}_{(at, dt)} \left[\prod_{l=1}^K \left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)}} \right)^{\alpha^{(l)}} \right] \right)^{N-S} \\ &\times \prod_{i=1}^S \mathbb{E}_{(at, dt)} \left[\prod_{l=1}^K \frac{\beta^{\alpha^{(l)}} \Gamma(\alpha^{(l)} + k_i^{(l)})}{\Gamma(\alpha^{(l)}) (\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)})^{\alpha^{(l)} + k_i^{(l)}}} \prod_{j=1}^M \frac{(t_{at, dt}^{j, (l)})^{k_{ij}^{(l)}}}{k_{ij}^{(l)}!} \right]. \end{aligned}$$

Note that in the general case where individuals' exact entrance/departure times are not known, we compute the above expectation over the entrance/departure time distribution $f(at, dt)$. In the special case where individuals' exact arrival/departure times are known, the likelihood function can be simplified as follows:

$$L = \binom{N}{N-S} \left(1 - \mathbb{E}_{(at, dt)} \left[\prod_{l=1}^K \left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)}} \right)^{\alpha^{(l)}} \right] \right)^{N-S} \\ \times \prod_{i=1}^S \left\{ \prod_{l=1}^K \frac{\beta^{\alpha^{(l)}} \Gamma(\alpha^{(l)} + k_i^{(l)})}{\Gamma(\alpha^{(l)}) \left(\beta + \sum_{j=1}^M t_{at, dt}^{j, (l)} \right)^{\alpha^{(l)} + k_i^{(l)}}} \right\}.$$

III. ADDITIONAL MEASUREMENTS

In this section we present additional measurements for the Paléo and the EPFL datasets.

A. Additional Measurements for the Paléo Dataset

Figure 3 shows the total number of different Bluetooth devices discovered, and the duration of stay (in minutes) on the festival grounds for each agent.

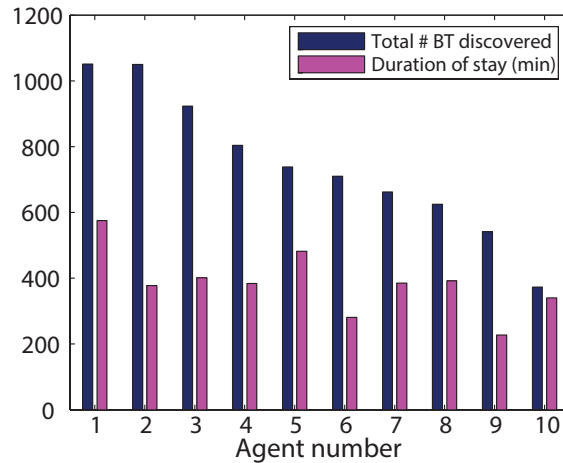


Fig. 3. The total number of different Bluetooth devices discovered, and the duration of stay (in minutes) on the festival grounds for each agent.

B. Additional Measurements for the EPFL Dataset

1) *Users' Arrival and Departure Times:* The empirical joint arrival/departure time distribution of the individuals (users) is plotted in Figure 4. The empirical marginal distributions of the individuals' arrival (first connection) times, departure (last connection) times, and the duration of stay on the campus is plotted in Figure 5.

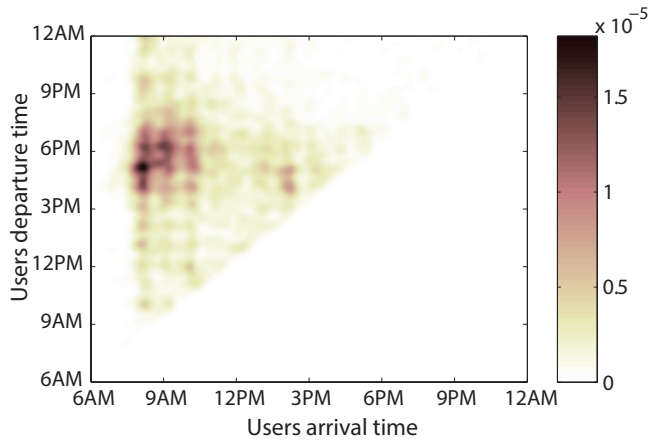


Fig. 4. Empirical distribution of individuals' joint arrival/departures times to/from the EPFL campus

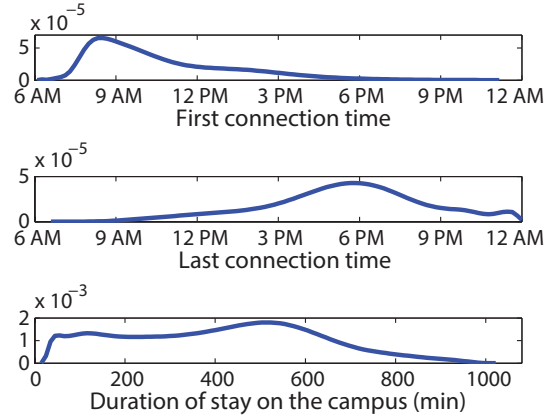


Fig. 5. Empirical marginal distributions of the individuals' arrival times (top), departure times (middle), and the duration of stay (bottom).

2) *Population Size Estimation using Subset of the Information:* Here, similar to Section 5.3 of [1], we apply our population size estimator to the EPFL dataset by varying the number of agents and the observation window. We first perform an experiment where we vary the number of agents from $M = 2$ to $M = 16$. For each value of M we simulate agents' trajectories as described in Section 8.3 of [1]; we assume that agents arrive at 06h00 and depart at 24h00. The average of estimated population sizes and their 90% confidence intervals, based on 1000 performed iterations, are shown in Figure 6 as a function of number of agents M . In another experiment we choose $M = 7$ agents but vary the duration of time that agents stay on the campus (the observation window). We consider the measurements obtained by all agents during an observation window of length w between 6 to 18 hours (the maximum possible duration is 18 hours). We partition the interval from 6h00 to 24h00 into 1-hour slots. For an observation window of length w , we consider all the consecutive 1-hour slots with total length w , and we estimate the population size based on the measurements obtained during those slots. For each window of length w we iterate over 1000 iteration, where at each iteration we simulate trajectories for $M = 7$ agents. The average of the estimated population sizes and their 90% confidence intervals for all the observation windows with length w are shown in Figure 7. The behaviors are close to those of the Paléo dataset plotted in Figure 5 of [1]: there exists an overshoot for small M . We have a very good estimate for the population size with a negligible bias for number of agents $M = 10$. Note that for the right-most point in both figures, there are 1000 estimates for the

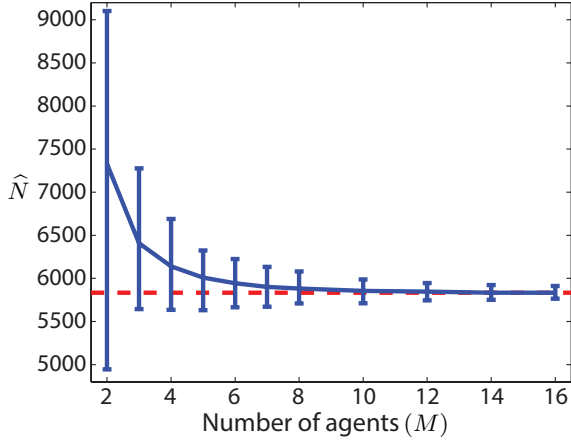


Fig. 6. Estimated population size as a function of the number of agents. The bars show the 90% confidence interval. The window length is 18 hours. The dashed line shows the number of agents is $M = 7$. The ground-truth for the population size.

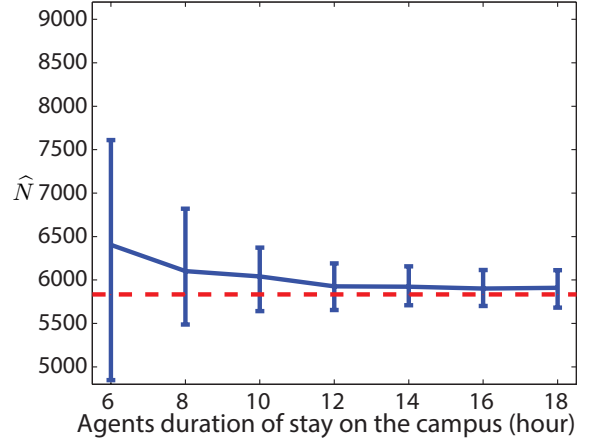


Fig. 7. Estimated population size as a function of the observation window. The bars show the 90% confidence interval. The observation window is 18 hours. The dashed line shows the number of agents is $M = 7$. The ground-truth for the population size.

population size, thus there is a non-zero standard deviation (in contrary to Figure 5 of [1]).

IV. PROOF OF THE THEOREMS

A. Background

In our proofs we use the following two theorems and proposition; for more information refer to textbooks in Statistics such as [2], [3], [4].

Theorem 1 (Fisher-Neyman Factorization Theorem). *Suppose that the observations (denoted by \mathbf{O}) have a joint density or frequency function $f(\mathbf{O}; \boldsymbol{\theta})$, where $\boldsymbol{\theta}$ is the vector of the parameters. A statistic $T = T(\mathbf{O})$ is sufficient for $\boldsymbol{\theta}$ if and only if the function f admits the factorization $f(\mathbf{O}; \boldsymbol{\theta}) = g(T(\mathbf{O}), \boldsymbol{\theta})h(\mathbf{O})$.*

Theorem 2 (Minimally sufficient Statistic). *Let the observations \mathbf{O} have joint density or frequency function $f(\mathbf{O}; \boldsymbol{\theta})$ and $T = T(\mathbf{O})$ be a statistic. Suppose that $f(\mathbf{O}; \boldsymbol{\theta})/f(\mathbf{O}'; \boldsymbol{\theta})$ is independent of $\boldsymbol{\theta}$ if and only if $T(\mathbf{O}) = T(\mathbf{O}')$. Then T is minimally sufficient for $\boldsymbol{\theta}$.*

B. Proofs

In the following, we prove the theorems in the paper.

Theorem 3. *The following quantities are the minimally sufficient statistics for estimating the population size in our model.*

- 1) *Number of times agent j contacts individual i (k_{ij}); note that only strictly positive values of k_{ij} 's are observed,*
- 2) *Agents' arrival/departure times to/from the festival,*
- 3) *Individuals' exact arrival/departure times to/from the festival, or their distribution $f(at, dt)$ or some approximation of the distribution.*

Proof. Recall that the full likelihood function had the following form,

$$L(\mathbf{O}; N, \alpha, \beta) = \binom{N}{S} \left(\mathbb{E}_{at, dt} \left[\left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at, dt}^j} \right)^\alpha \right] \right)^{N-S} \\ \times \prod_{i=1}^S \left\{ \frac{\Gamma(\alpha + k_i) \beta^\alpha}{\Gamma(\alpha)} \mathbb{E}_{at, dt} \left[\frac{\prod_{j=1}^M \frac{(t_{at, dt}^j)^{k_{ij}}}{k_{ij}!}}{(\beta + \sum_{j=1}^M t_{at, dt}^j)^{\alpha + k_i}} \right] \right\},$$

where \mathbf{O} represents the observed measurements. We first prove the sufficiency part by factorizing the likelihood function as follows,

$$L(\mathbf{O}; N, \alpha, \beta) = g(\mathbf{O}; N, \alpha, \beta) \cdot h(\mathbf{O}), \quad (2)$$

where

$$g(\mathbf{O}; N, \alpha, \beta) = \frac{N!}{(N-S)!} \left(\mathbb{E}_{at, dt} \left[\left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at, dt}^j} \right)^\alpha \right] \right)^{N-S} \\ \times \prod_{i=1}^S \left\{ \frac{\Gamma(\alpha + k_i) \beta^\alpha}{\Gamma(\alpha)} \mathbb{E}_{at, dt} \left[\frac{\prod_{j=1}^M (t_{at, dt}^j)^{k_{ij}}}{(\beta + \sum_{j=1}^M t_{at, dt}^j)^{\alpha + k_i}} \right] \right\},$$

and

$$h(\mathbf{O}) = \frac{1}{S!} \prod_{i=1}^S \left\{ \prod_{j=1}^M \frac{1}{k_{ij}!} \right\}.$$

The function $h(\mathbf{O})$ is only a function of the observations, and $g(\mathbf{O}; N, \alpha, \beta)$ is function of the parameters (α , β , and N) and the statistics defined in the theorem (k_{ij} 's, $f(at, dt)$, and agents' arrival/departure times). Note that k_i and S are functions of the statistics defined in the theorem (i.e., functions of k_{ij} 's). Thus the sufficiency result follows from the Fisher-Neyman Factorization Theorem based on the factorization in (2).

We now prove the minimally sufficiency part. Assume that we have two sets of experiments, for which the number of agents M and agents' respective arrival/departure times are the same.

Furthermore, assume that the individuals' arrival/departure time distribution $f(at, dt)$ (or its approximation) is also the same for both experiments. Denote the measurements for the first experiment by S, k_i for $i = 1, 2, \dots, S$, and k_{ij} , for $i = 1, 2, \dots, S, j = 1, 2, \dots, M$; and for the second experiment by S', k'_i for $i = 1, 2, \dots, S'$, and k'_{ij} , for $i = 1, 2, \dots, S', j = 1, 2, \dots, M$.

The minimally sufficiency of the statistics follows from Theorem 2 since the ratio

$$\begin{aligned} \frac{L(\mathbf{O}; N, \alpha, \beta)}{L(\mathbf{O}'; N, \alpha, \beta)} &= \frac{S!(N-S)!}{S!(N-S)!} \left(\mathbb{E}_{at,dt} \left[\left(\frac{\beta}{\beta + \sum_{j=1}^M t_{at,dt}^j} \right)^\alpha \right] \right)^{S'-S} \left(\frac{\beta^\alpha}{\Gamma(\alpha)} \right)^{(S-S')} \\ &\times \frac{\prod_{i=1}^S \left\{ \Gamma(\alpha + k_i) \mathbb{E}_{at,dt} \left[\frac{\prod_{j=1}^M (t_{at,dt}^j)^{k_{ij}}}{(\beta + \sum_{j=1}^M t_{at,dt}^j)^{\alpha+k_i}} \right] \right\}}{\prod_{i=1}^{S'} \left\{ \Gamma(\alpha + k'_i) \mathbb{E}_{at,dt} \left[\frac{\prod_{j=1}^M (t_{at,dt}^j)^{k'_{ij}}}{(\beta + \sum_{j=1}^M t_{at,dt}^j)^{\alpha+k'_i}} \right] \right\}} \cdot \frac{\prod_{i=1}^S \left\{ \prod_{j=1}^M \frac{1}{k_{ij}!} \right\}}{\prod_{i=1}^{S'} \left\{ \prod_{j=1}^M \frac{1}{k'_{ij}!} \right\}} \end{aligned}$$

is independent of N, α, β if and only if the statistics given in the theorem have the same set of values, i.e., when

$$S = S', \{k_i | i = 1, \dots, S\} = \{k'_i | i = 1, \dots, S'\},$$

and

$$\{[k_{i1}, k_{i2}, \dots, k_{iM}] | i = 1, \dots, S\} = \{[k'_{i1}, k'_{i2}, \dots, k'_{iM}] | i = 1, \dots, S'\}.$$

To see this, consider for example the case where $S \neq S'$. In this case, the above ratio will be a function of $\left(\frac{\beta^\alpha}{\Gamma(\alpha)}\right)^{(S-S')}$, which makes the ratio a function of the parameters α and β . Now if $S = S'$, we have

$$\frac{L(\mathbf{O}; N, \alpha, \beta)}{L(\mathbf{O}'; N, \alpha, \beta)} \propto \frac{\prod_{i=1}^S \{\Gamma(\alpha + k_i)\}}{\prod_{i=1}^{S'} \{\Gamma(\alpha + k'_i)\}} \cdot \frac{\prod_{i=1}^S \left\{ \mathbb{E}_{at,dt} \left[\frac{\prod_{j=1}^M (t_{at,dt}^j)^{k_{ij}}}{(\beta + \sum_{j=1}^M t_{at,dt}^j)^{\alpha+k_i}} \right] \right\}}{\prod_{i=1}^{S'} \left\{ \mathbb{E}_{at,dt} \left[\frac{\prod_{j=1}^M (t_{at,dt}^j)^{k'_{ij}}}{(\beta + \sum_{j=1}^M t_{at,dt}^j)^{\alpha+k'_i}} \right] \right\}}.$$

□

Now if $\{k_i | i = 1, \dots, S\} \neq \{k'_i | i = 1, \dots, S'\}$, then the first term depends on parameter α . We can similarly argue that if $S = S'$ and $\{k_i | i = 1, \dots, S\} = \{k'_i | i = 1, \dots, S'\}$ but $\{[k_{i1}, k_{i2}, \dots, k_{iM}] | i = 1, \dots, S\} \neq \{[k'_{i1}, k'_{i2}, \dots, k'_{iM}] | i = 1, \dots, S'\}$, again the above ratio will be a function of α and β .

Theorem 4. *The following quantities are the minimally sufficient statistics for jointly estimating the population size and density in our model.*

- 1) *Number of times agent j contacts individual i in location l ; note that only strictly positive values of $k_{ij}^{(l)}$'s are observed,*
- 2) *Agents' arrival/departure times to/from the festival, and their trajectories on the festival grounds,*
- 3) *Individuals' exact arrival/departure times to/from the festival, or their distribution $f(at, dt)$ or some approximation of the distribution.*

Proof. Proof is exactly similar to the proof of Theorem 3. □

REFERENCES

- [1] F. Movahedi Naini, O. Dousse, P. Thiran, and M. Vetterli, "Opportunistic sampling for joint population size and density estimation," *IEEE Transactions on Mobile Computing*, 2014.
- [2] K. Knight, *Mathematical Statistics*, ser. Texts in Statistical Science Series. Taylor & Francis, 2010. [Online]. Available: <http://books.google.ch/books?id=VSfdpgTZScwC>
- [3] P. Bickel and K. Doksum, *Mathematical Statistics: Basic Ideas And Selected Topics*, ser. Mathematical Statistics: Basic Ideas and Selected Topics. Prentice Hall, 2006, no. v. 1. [Online]. Available: <http://books.google.ch/books?id=wImDQgAACAAJ>
- [4] D. Cox and D. Hinkley, *Theoretical Statistics*. Chapman and Hall, 1979. [Online]. Available: <http://books.google.ch/books?id=ppoujo-BInsC>