



Sampling Size Optimization for Bioburden Density Estimation in Planetary Protection

July 2024

Changing the World's Energy Future

Lisa Guan, Andrei Vasilyevich Gribok, Michael DiNicola



DISCLAIMER

This information was prepared as an account of work sponsored by an agency of the U.S. Government. Neither the U.S. Government nor any agency thereof, nor any of their employees, makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness, of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. References herein to any specific commercial product, process, or service by trade name, trade mark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the U.S. Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the U.S. Government or any agency thereof.

Sampling Size Optimization for Bioburden Density Estimation in Planetary Protection

Lisa Guan, Andrei Vasilyevich Gribok, Michael DiNicola

July 2024

**Idaho National Laboratory
Idaho Falls, Idaho 83415**

<http://www.inl.gov>

**Prepared for the
U.S. Department of Energy
Under DOE Idaho Operations Office
Contract DE-AC07-05ID14517**

Jet Propulsion Laboratory
California Institute of Technology



Sampling Size Optimization for Bioburden Density Estimation in Planetary Protection

Andrei Gribok^a, Mike DiNicola^b, and Lisa Guan^b

^a Idaho National Laboratory, Idaho Falls, ID, USA

^b Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA

45th Scientific Assembly of the Committee on Space Research (COSPAR), July
13 - July 21, 2024, BEXCO, Busan, Republic of Korea.

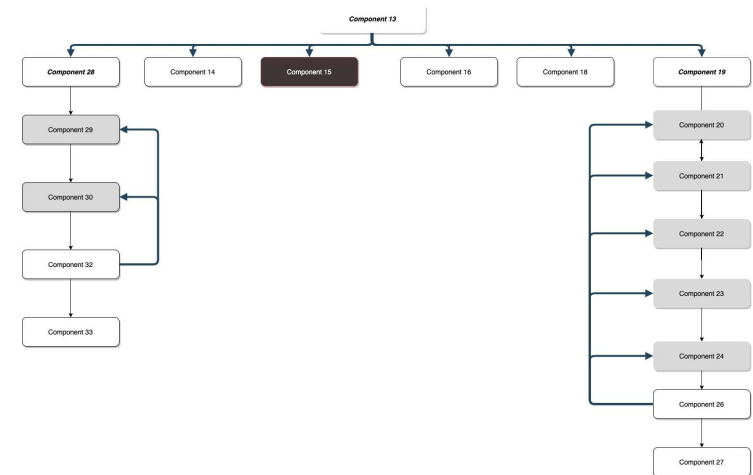
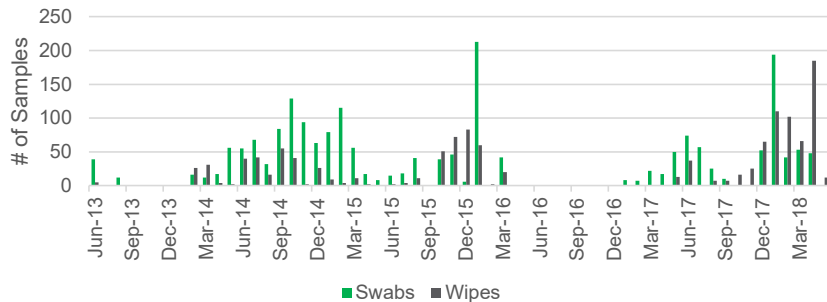
Battelle Energy Alliance manages INL for the
U.S. Department of Energy's Office of Nuclear Energy



Idaho National Laboratory

Planetary Protection Dataset Overview – Interior Exploration using Seismic Investigations, Geodesy and Heat Transport (InSight) Mission

- ~10% of spacecraft surface areas were verified for biological cleanliness.
 - 2,031 swabs, 1,266 wipes, and 39,379 petri dishes were processed between June 2013 and May 2018.
 - 118,137 data points were recorded as part of the biological cleanliness verification campaign.
- A majority of the samples yielded clean results.
 - 93% of the swabs and 63% of the wipes had 0 colony forming units (CFUs), resulting in ~85% of the 39,379 petri dishes yielding 0 CFUs (Hendrickson et al., Astrobiology 2019).



Tree structure of hardware component 13. The sampled subcomponents are shown as white rectangles, implied components are denoted in grey, and the specified component is shown in black. The names of rollup components 13, 19, and 28 are italicized on the diagram.

Planetary Protection Biological Cleanliness Verification Sampling



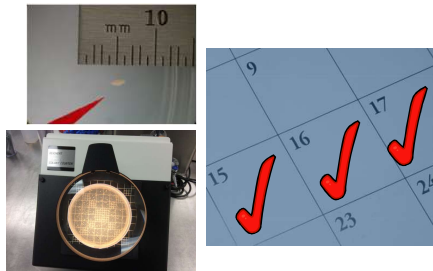
Clean hardware (and table or bag) with a solvent wipe prior to sampling or installation



Sample with swab or wipe (water is used as the solvent)



Process swabs and wipes, up to 24 hours required post-assay



Count plates – 24 h, 48 h, and 72 h

- Swabs sampled a 0.0025 m² surface area maximum.
- Wipes sampled up to a 1.0 m² surface area.
- Due to the experimental procedure, the swabs assume a pour fraction of 0.8 and the wipes 0.25, representing the portion of the total sample solution plated and analyzed for CFU counts.
- Exposure is the sampled area multiplied by the pour fraction (i.e., for swabs it would be $E = 0.0025 \text{ m}^2 \cdot 0.8 = 0.002 \text{ m}^2$).

Spore/CFU, a heat-tolerant reproductive cell capable of developing into a new individual without fusion with another reproductive cell.

Data Generating Model and Parameter Estimators

$$X_i = x | \lambda_{true}^i \sim \text{Poisson}(\lambda_{true}^i \cdot E_i)$$

$$\lambda_{true}^i | \alpha, \beta \sim \text{Gamma}(\alpha, \beta), i = 1, \dots, N$$

$$\hat{\lambda}_i(x_i) = \frac{x_i}{E_i} - \text{Maximum Likelihood Estimator (MLE)}$$

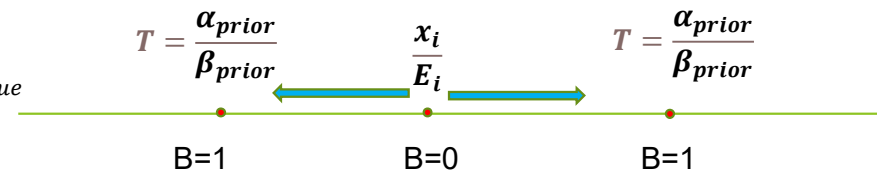
$$\hat{\lambda}_i(x_i) = \frac{x_i + \alpha}{E_i + \beta} - \text{Conjugate prior Bayes Estimator}$$

$$\hat{\lambda}_i = d - \text{Deterministic Estimator}$$

- Bayesian and Deterministic estimators are shrinkage estimators with respect to MLE
- Shrinkage estimators pulls MLE towards a prespecified target value
- Shrinkage reduces the variance of the estimate at the expense of introducing bias
- Variance reduction can outweigh increase in bias thus improving mean squared error of the estimator

$$\hat{\lambda}_i = \frac{x_i}{E_i} - B \cdot \left(\frac{x_i}{E_i} - T \right) \quad B = \frac{\beta_{prior}}{E_i + \beta_{prior}} \leq 1$$

X_i – random variable describing the number of CFUs on i-th sampled component
 x_i – observed number of CFUs in a sample from i-th sampled component
 E_i – exposure for a sample, determined by multiplying the sampled area by pour fraction
 λ_{true}^i – true bioburden density for a sample
 $\text{Gamma}(\lambda | \alpha, \beta)$ – Gamma prior distribution of bioburden density λ_{true}^i
 α, β – parameters of the Gamma prior distribution
 $\text{Poisson}(\lambda_{true}^i \cdot E_i)$ – Poisson likelihood of having x CFU counts at exposure E_i given λ_{true}^i
 N – number of components



Estimators' Loss Functions and Risks

$$L(\hat{\lambda}_i(x), \lambda_{true}^i) = (\hat{\lambda}_i(x) - \lambda_{true}^i)^2 - \text{Loss function for a data - driven estimators, (MLE and Bayes)}$$

$$L(\hat{\lambda}_i, \lambda_{true}^i) = (\hat{\lambda}_i - \lambda_{true}^i)^2 - \text{Loss function for a deterministic estimator}$$

$$R_{Frequentist}(\lambda_{true}^i, \hat{\lambda}_i(x)) = \sum_x (\lambda_{true}^i - \hat{\lambda}_i(x))^2 \cdot p(x; \lambda_{true}^i) = \sum_{x=0}^{\infty} (\lambda_{true}^i - \hat{\lambda}_i(x))^2 \cdot \frac{(\lambda_{true}^i \cdot E_i)^x}{x!} \cdot e^{-\lambda_{true}^i \cdot E_i}$$

$$\rho_{Posterior \text{ expected loss}}(\alpha, \beta, x_i, \hat{\lambda}_i(x)) = \int_0^{\infty} (\lambda_{true}^i - \hat{\lambda}_i(x_i))^2 \cdot \text{Gamma}(\lambda_{true}^i | x_i, \alpha, \beta) d\lambda_{true}^i$$

$$r_{Integrated}(\alpha, \beta, \hat{\lambda}_i) = \int_0^{\infty} R(\lambda_{true}^i, \hat{\lambda}_i(x)) \cdot \text{Gamma}(\lambda_{true}^i | x_i, \alpha, \beta) d\lambda_{true}^i = \sum_{x=0}^{\infty} \rho(\alpha, \beta, x_i, \hat{\lambda}_i(x)) \cdot NB(x | \alpha, \frac{\beta}{\beta + E_i})$$

Risks for Three Different Estimators

Estimator	R-Frequentist Risk (FR)	ρ -Posterior Expected Loss (PEL)	r- Integrated Risk (IR)
d (deterministic)	$(d - \lambda_{true})^2$	$\frac{\alpha}{\beta^2} + \left(d - \frac{\alpha}{\beta}\right)^2$	$\frac{\alpha}{\beta^2} + \left(d - \frac{\alpha}{\beta}\right)^2$
$\frac{x}{E}$ (MLE)	$\frac{\lambda_{true}}{E}$	$\frac{x + \alpha}{(E + \beta)^2} + \left(\frac{x}{E} - \frac{\alpha}{\beta}\right)^2$	$\frac{\alpha}{E \cdot \beta}$
$\frac{x + \alpha}{E + \beta}$ (Bayes)	$\frac{\lambda_{true}}{E} \cdot (1 - B)^2 + \left(B \cdot \left(\lambda_{true} - \frac{\alpha}{\beta}\right)\right)^2$ $B = \frac{\beta}{E + \beta} \leq 1$	$\frac{x + \alpha}{(E + \beta)^2}$	$\frac{\alpha}{\beta \cdot (E + \beta)}$
Properties of the Risks	<ul style="list-style-type: none"> The FR is a function of unknown parameter λ_{true} Functions are harder to compare than numbers Still useful, for example, $\lambda_{true} < 2 \cdot \lambda_{true}$ 	<ul style="list-style-type: none"> PEL depends only on known values (data x) PEL is a number PEL may not exist for some types of prior 	<ul style="list-style-type: none"> IR is a number IR depends neither on data nor λ_{true} IR can be used for sampling optimization

* From a dimensional analysis perspective, risk and loss functions above have units [CFU/Exposure]². This can be seen to be consistent with the above formulas via appropriate transformation; e.g., for the R-Frequentist Risk of the MLE estimator, $\frac{\lambda_{true}}{E} = \frac{\lambda_{true}E}{E^2}$.

Integrated Risk and Sampling Size Determination for Implied Components

$$TR_{Bayes}(\alpha, \beta, E) = \underbrace{R(\alpha, \beta, E)_{Bayes}}_{\text{Uncertainty}} + \underbrace{C_0 + E \cdot C}_{\text{Cost}} = \frac{k \cdot \alpha}{\beta \cdot (E + \beta)} + C_0 + C \cdot E$$

For swabs used for the InSight mission, the cost C was \$9.11 for the first swab and \$1.45 for each additional swab
 C₀-sampling setup cost
 k-scaling factor that dollarizes the r-integrated risk

$$E_{opt} = \max \left\{ 0, \sqrt{\frac{k \cdot \alpha}{C \cdot \beta}} - \beta \right\}$$

Implied Component #	Implied from Component # (implicant)	Implied Bioburden Density, $\hat{\lambda}$, CFUs/m ²	$\sqrt{\text{Implied Risk}}$, CFUs/m ²	Total Surface Area of the Implied Component, m ²	Optimal Sampling Area, m ²	$\sqrt{\text{Optimal Risk}}$, CFUs/m ²	Optimal Cost, \$
2	10	12.05	17.03	0.26	0.16	7.68	126.80
106	108	13.51	19.11	0.53	0.17	7.91	138.78
133	131	5.10	4.16	0.013	0	4.16	0
36	38	15.50	4.38	2.0	0	4.38	0
71	70	2.47	2.02	7.0	0	2.02	0
29	32	104.16	65.88	0.16	0.16	23.41	129.46

Analysis of InSight's Sampled Components

Component	Bioburden density, $\hat{\lambda}$, CFUs/m ²	\sqrt{Risk} , CFUs/m ²	Total surface area m ²	Area sampled, m ²	Optimal additional sampling area, m ²	$\sqrt{Optimal Risk}$, CFUs/m ²	Optimal Cost, Units
Component 233	31.25	44.19	0.0285	0.02	0.0285	26.49	2.44
Component 236	173.07	81.58	0.5	0.0325	0.17	29.48	13.86
Component 237	93.75	76.54	0.5	0.02	0.14	24.23	11.65
Component 238	125.00	176.77	0.0179	0.005	0.0179	75.54	1.54
Component 242	15.62	22.09	0.298	0.04	0.051	13.71	4.32
Component 243	48.24	20.57	0.298	0.28	0.003	20.29	0.27
Component 245	35.71	50.50	0.282	0.0175	0.098	17.80	8.16
Component 246	62.50	88.38	0.237	0.01	0.130	21.25	10.63
Component 26	27.55	14.73	0.454	0.255	6.61E-05	14.72	0.0057
Component 46	83.33	117.85	0.024	0.007	0.024	52.70	2.05
Component 49	46.87	38.27	0.26	0.04	0.0907	19.54	7.54

Summary

- The quality of an estimator is measured by its risks
- The risk can be defined as a mathematical expectation of the difference between the estimate and the true value of the parameter
- Three different risks are commonly used in statistics: frequentist, Bayesian , and integrated
- The integrated risk depends on exposure and parameters of prior distribution
- The integrated risk is a valuable tool to optimize sampling for implied, specified, and sampled components
- The future work will include multi-dimensional risk-cost optimization for multiple components to obtain optimal bioburden estimation strategies

Acknowledgements



This work was supported by U.S. DOE-NASA Strategic Partnership Project (SPP) #19701.

The authors are grateful to Dr. J. Nick Benardini, Dr. Elaine E. Seasley for their contribution to this project.