



**HAL**  
open science

# A computational framework to systematize uncertainty analysis in the sediment fingerprinting approach using least square methods

Lidiane Buligon, Tiago Martinuzzi Buriol, Jean Paolo Gomes Minella, Olivier Evrard

## ► To cite this version:

Lidiane Buligon, Tiago Martinuzzi Buriol, Jean Paolo Gomes Minella, Olivier Evrard. A computational framework to systematize uncertainty analysis in the sediment fingerprinting approach using least square methods. *Journal of Computational and Applied Mathematics*, 2024, 43 (8), pp.444. 10.1007/s40314-024-02948-4 . cea-04748250

**HAL Id: cea-04748250**

**<https://cea.hal.science/cea-04748250v1>**

Submitted on 28 Oct 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Computational and Applied Mathematics

## A computational framework to systematize uncertainty analysis in the sediment fingerprinting approach using Least Square Methods --Manuscript Draft--

<b>Manuscript Number:</b>							
<b>Full Title:</b>	A computational framework to systematize uncertainty analysis in the sediment fingerprinting approach using Least Square Methods						
<b>Article Type:</b>	Original Article						
<b>Funding Information:</b>	<table border="1"><tr><td>CEA-Atomic Energy Commission</td><td>Mr Olivier Evrard</td></tr><tr><td>National Council for Scientific and Technological Development (200008/2023-4)</td><td>Mrs Lidiane Buligon</td></tr><tr><td>CAPES-Coordination for the Improvement of Higher Education Personnel (88887.310201/2018-00)</td><td>Mrs Lidiane Buligon</td></tr></table>	CEA-Atomic Energy Commission	Mr Olivier Evrard	National Council for Scientific and Technological Development (200008/2023-4)	Mrs Lidiane Buligon	CAPES-Coordination for the Improvement of Higher Education Personnel (88887.310201/2018-00)	Mrs Lidiane Buligon
CEA-Atomic Energy Commission	Mr Olivier Evrard						
National Council for Scientific and Technological Development (200008/2023-4)	Mrs Lidiane Buligon						
CAPES-Coordination for the Improvement of Higher Education Personnel (88887.310201/2018-00)	Mrs Lidiane Buligon						
<b>Abstract:</b>	<p>Simulating sediment transfer processes in catchments has contributed significantly to solving environmental problems due to its importance in the silting of rivers and reservoirs and for controlling the pollution of water bodies. Among the methods used to improve data collection and modelling, the "sediment fingerprinting approach" uses tracers reflecting the composition of eroded soils and sediments in multivariate statistical analyses and mathematical models for optimizing equation systems. Based on generalized least squares (GLS) method and Mahalanobis distance, this study sought to present a computational framework to solve over-determined systems applied to sediment tracing, systematize the uncertainty analysis and sample number optimization. Hence, this approach takes into account the influence of collinearity among the chemical variables that compose the tracer set to be evaluated by the presence of the variance-covariance matrix. A dataset from the Arvorezinha experimental catchment in southern Brazil was used to validate the modeling, and our findings confirmed the assumption of increased uncertainty as the number of target samples decreases in the sources or eroded sediment samples. Sharing the file with the Python code contributes to improving the technique as it allows other researchers to systematically improve the definition of the number of samples required based on the uncertainty analysis.</p>						
<b>Corresponding Author:</b>	Lidiane Buligon, Dr Federal University of Santa Maria: Universidade Federal de Santa Maria Santa Maria, Rio Grande do Sul BRAZIL						
<b>Corresponding Author Secondary Information:</b>							
<b>Corresponding Author's Institution:</b>	Federal University of Santa Maria: Universidade Federal de Santa Maria						
<b>Corresponding Author's Secondary Institution:</b>							
<b>First Author:</b>	Lidiane Buligon, Dr						
<b>First Author Secondary Information:</b>							
<b>Order of Authors:</b>	Lidiane Buligon, Dr Tiago Martinuzzi Buriol, Dr Jean Paolo Gomes Minella, Dr Olivier Evrard, PhD						
<b>Order of Authors Secondary Information:</b>							
<b>Author Comments:</b>							
<b>Suggested Reviewers:</b>	Gabriel Haeser ghaeser@ime.usp.br						

Research Interests: Nonlinear Optimization, Conic Optimization, Numerical Analysis, Optimality Conditions, Algorithms, Applications.

Patrick Laceby  
patrick.laceby@gov.ab.ca

Skills and expertise:  
Sediments, Geochemical Modeling Stable Isotopes, Sedimentology, Geochemistry  
Stable, Isotope Analysis

Arman Haddadchi  
aman.haddadchi@niwa.co.nz

Skills and expertise:  
Genetic Algorithm, Soil Erosion ,Sediment, Sediments River, Engineering Environment  
,Geomorphology

Leonardo Ramos Emmendorfer  
leonardo.emmendorfer@gmail.com

Skills and expertise: Evolutionary Computation, Clustering, Evolutionary Algorithms,  
Heuristics, Optimization

[Click here to view linked References](#)

1  
2  
3  
4  
5  
6 1 A computational framework to systematize  
7  
8 2 uncertainty analysis in the sediment  
9  
10 3 fingerprinting approach using Least Square  
11  
12 4 Methods

13  
14  
15 5 Lidiane Buligon<sup>1\*</sup>, Tiago Martinuzzi Buriol<sup>1</sup>,  
16 6 Jean Paolo Gomes Minella<sup>2</sup>, Olivier Evrard<sup>3†</sup>

17  
18 7 <sup>1</sup>Department of Mathematics, Federal University of Santa Maria,  
19 8 Roraima Av., n.1000, Santa Maria, 97105-900, Rio Grande do Sul, Brazil.

20 9 <sup>2</sup>Department of Soils, Federal University of Santa Maria, Roraima Av.,  
21 10 n.1000, Santa Maria, 97105-900, Rio Grande do Sul, Brazil.

22 11 <sup>3</sup>Laboratoire des Sciences du Climat et de l'Environnement (LSCE-IPSL),  
23 12 CEA Saclay, Orme des Merisiers, Gif-sur-Yvette, 91 191 Cedex, France.

24  
25  
26  
27 13 \*Corresponding author(s). E-mail(s): [buligon.l@ufsm.br](mailto:buligon.l@ufsm.br);  
28 14 Contributing authors: [tiago.buriol@ufsm.br](mailto:tiago.buriol@ufsm.br); [jean.minella@ufsm.br](mailto:jean.minella@ufsm.br);  
29 15 [olivier.evrard@lsce.ipsl.fr](mailto:olivier.evrard@lsce.ipsl.fr);

30 16 † These authors contributed jointly to this work

31  
32  
33  
34 17 **Acknowledgments**

35 18 The authors acknowledge the financial support granted by the CNPq (National  
36 19 Council for Scientific and Technological Development) and CAPES (Coordina-  
37 20 tion for the Improvement of Higher Education Personnel) by the first author's  
38 21 scholarship received during the development of this study. CEA (Atomic Energy  
39 22 Commission) in supporting a 1-year research visit in France. We would also like  
40 23 to thank Atlas Assessoria Linguística for support with the English version of this  
41 24 manuscript.

42  
43  
44 25 **Abstract**

45 26 Simulating sediment transfer processes in catchments has contributed signifi-  
46 27 cantly to solving environmental problems due to its importance in the silting of  
47 28 rivers and reservoirs and for controlling the pollution of water bodies. Among

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

the methods used to improve data collection and modelling, the “sediment fingerprinting approach” uses tracers reflecting the composition of eroded soils and sediments in multivariate statistical analyses and mathematical models for optimizing equation systems. Based on generalized least squares (GLS) method and Mahalanobis distance, this study sought to present a computational framework to solve over-determined systems applied to sediment tracing, systematize the uncertainty analysis and sample number optimization. Hence, this approach takes into account the influence of collinearity among the chemical variables that compose the tracer set to be evaluated by the presence of the variance-covariance matrix. A dataset from the Arvorezinha experimental catchment in southern Brazil was used to validate the modeling, and our findings confirmed the assumption of increased uncertainty as the number of target samples decreases in the sources or eroded sediment samples. Sharing the file with the Python code contributes to improving the technique as it allows other researchers to systematically improve the definition of the number of samples required based on the uncertainty analysis.

**Keywords:** Computational Mathematics, Generalized Least Squares, Sediment Source Identification, Mahalanobis Distance, Confidence Region

## 1 Introduction

Water erosion is one of the main processes of soil degradation, severely impacting water resources [1]. Researchers have demonstrated that modelling sediment production and transfer requires a set of mathematical techniques that rely on robust monitoring network [2]. The modelling techniques depend on the efficiency of the measured data in representing the spatial variability of sediment sources. In addition, the models seek to represent the temporal variability of eroded sediment mobilization and transfer phenomena. Thus, the ability of the models to represent the spatial and temporal variability of the phenomena depends essentially on the data sets available. On the other hand, the model outputs should maximize the explanatory capacity of the sample set, minimizing the costs associated with data collection and analysis. The sediments in the river originate from various locations where the erosion process can occur with specific patterns and magnitudes [3]. Therefore, quantitatively defining the contribution of each source is crucial to propose effective soil conservation measures and reduce environmental and economic problems linked to soil erosion.

Sediment source identification is an important branch of sediment production modeling that employs tracers found in soils and sediments in a set of statistical analysis techniques known as the “sediment source identification” or “sediment fingerprinting/tracing approach” [4–10]. This approach involves various areas of knowledge, including geoscience, statistics, and computational mathematics, fostering the exchange of knowledge from the choice of potential tracers to the use of mathematical and statistical methods to describe the different steps of the process. Optimizing models based on observing the processes improves the estimation of contributions from different sediment sources, quantifying the associated uncertainties and extracting the

70 maximum amount of information from the data set obtained in the field. Further-  
71 more, the shared use of open-source algorithms expands the applications of these new  
72 methods among different research groups, thereby enhancing the scientific and tech-  
73 nological advances in soil science and related disciplines. Numerous contributions have  
74 been made to improving techniques for identifying sediment sources for monitoring  
75 and modelling [3, 11–17]. The main assumptions of the suspended sediment identifi-  
76 cation technique are that: (1) sources can be discriminated by different characteristics  
77 (tracer properties) found in the source soils; (2) the eroded sediments in the river  
78 system consist of a mixture of the sediments originating from potential sources accu-  
79 mulated during the transfer process across the catchment; (3) the temporal variations  
80 in the tracer properties of the eroded sediments found in the river directly reflect the  
81 spatial variation of the erosion processes according to the behavior of each source; (4)  
82 the tracer properties of the sources and suspended sediments can be compared to es-  
83 tablish the contribution of each source to the sampled sediment (target). Given these  
84 assumptions, the sediment fingerprinting approach establishes a relationship between  
85 source characteristics and those of suspended sediments by mathematically solving an  
86 over-determined linear system of equations (i.e., linear systems in which the number  
87 of equations is greater than the number of unknowns).

88 Applications in different areas of research are described by over-determined linear  
89 systems, including the traditional modeling of sediment source identification proposed  
90 by Walling and Woodward [18]. Methods for solving this type of system of equations  
91 have been continuously improved in order to obtain better estimates for the set of ob-  
92 served quantities. It is expected that the more information available, the better the  
93 quality of the results obtained from your analysis. However, obtaining feasible solu-  
94 tions for over-determined systems is a challenge. Regression analysis has been the most  
95 used technique in these cases, as it measures the direction and intensity of the relation-  
96 ship between the dependent and independent variables and numerically describes this  
97 relationship [19–22]. The methods use the mean square error to calculate and evaluate  
98 the performance of an estimator. In the ordinary least squares (OLS) method, a sys-  
99 tem of linear equations corresponds to a matrix equation of the form  $Ax = b$ , where  
100 the matrix  $A$  and the vector  $b$  are given, and  $x$  is the unknown solution [23–25]. The  
101 OLS algorithm was employed to create the FingerPro package [26], which was devel-  
102 oped in the R language to determine the contribution of sediment sources to target  
103 material. However, when there is a certain degree of correlation between the residu-  
104 als in a regression model, the Aitken estimator (or Gauss-Markov estimator) should  
105 be applied. In practice, however, the covariance matrix of the error is generally un-  
106 known, making this estimator nonviable. In such cases, a Generalized Least Squares  
107 (GLS) estimator, which is defined as the Gauss-Markov estimator with the unknown  
108 covariance matrix replaced with a suitable estimator, is used [21]. Least square meth-  
109 ods can lead to the obtainment of approximate solutions of over-determined systems,  
110 although sometimes no exact solution can be found. The basis for obtaining solutions  
111 to these types of systems is based on matrix algebra, which computational advances  
112 have been greatly boosted during the last several decades [27–29].

113 According to Walling [4] one of the challenges in improving the sediment source  
114 identification technique is related to the estimation of the uncertainties of the results.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

115 From this perspective, the influence of the number of samples used to characterize both  
116 the potential sources and the target sediment is critical for evaluating the confidence  
117 level associated with the results. The assumption, in this case, is based on estimat-  
118 ing the increase in uncertainty when the number of samples decreases. In addition,  
119 the number of tracers selected and their relationship with the number of sources are  
120 also topics under investigation [2]. Latorre et al. [13] emphasized that employing more  
121 tracers than the number of potential sources can lead to mathematical inconsistencies  
122 due to obtaining multiple solutions in over-determined systems. Although solutions of  
123 over-determined systems provide a possible source of errors for the sediment finger-  
124 printing approach, the authors believe that this topic has remained under-investigated  
125 based on the premise that the selection of tracers and the potential advantage pro-  
126 vided by the use of different models (e.g., FingerPro [26], MixSIAR [30], and SIFT  
127 [31]) have led to the primary advances of the technique.

128 In other studies, Haddadchi et al. [32] and Haddadchi et al. [33] compared the ac-  
129 curacy of different mixing models relying on the same source and target sediment data  
130 set. The results indicated that the relative contribution of sources to sediments de-  
131 pend on model was applied. In turn, Latorre et al. [13] and Lizaga et al. [12] pointed  
132 out that the use of different tracer selection methods may affect equally the outputs  
133 of the models. However, when using the same selection process, both types of models  
134 produced similar output results. Both studies agreed that the models based on Finger-  
135 Pro and MixSIAR did not test for consistency and the existence of multiple solutions  
136 in over-determined systems; for this reason, they recommend the implementation of  
137 new tracer selection method. In this context, Latorre et al. [13] designed the Consis-  
138 tent Tracer Selection (CTS) method to extract the solutions in the datasets of each  
139 mixture (from the selection of subsets of tracers). According to Latorre et al. [13], the  
140 outputs of CTS method were efficient in extracting solutions for the study case.

141 The methods applied to quantify the uncertainty associated with sediment fin-  
142 gerprinting are commonly based on the Monte Carlo simulation technique. For each  
143 tracer from each source groups/sink an estimate of the mean is obtained considering  
144 its probability distributions, and thus optimized solutions of the mixing model are  
145 obtained repeatedly [6, 32, 34–36]. From this modelling Franks and Rowan [6] devel-  
146 oped a model that allowed to determine the confidence intervals on un-mixing model  
147 results. In this analysis, the variability of each source and number of samples used in-  
148 fluenced the model performance. However, according to [37] it is necessary to consider  
149 weightings to take into account of the within-source variability and discriminatory  
150 power of individual tracer properties and correction factors (particle size and organic  
151 matter content for source type or spatial source category) to calculate the relative  
152 contribution of sources. In another study, Laceby and Olley [5] proposed a distri-  
153 bution mixing model, whose approach reincorporates correlations between elemental  
154 concentrations and models distributions for source contribution terms for multiple  
155 targets (end-members). This study demonstrated how different weightings can affect  
156 modelling results. Nevertheless, correlation between tracers are not considered.

157 In this sense, Clarke [36] and Clarke and Minella [38] proposed a statistical method  
158 that quantifies the number of samples required to characterize sources and sediments  
159 on estimating the uncertainties of the results. In addition, they added the effect of

160 collinearity between tracer variables to the analysis. In this approach, the variance-  
161 covariance matrix is used to solve the over-determined system; the authors included  
162 the correlations between variables (covariance) and the effect of the variance in each  
163 tracer variable. They also used the Mahalanobis distance to determine the confidence  
164 region associated with the uncertainties, and the model shows how the area of the con-  
165 fidence region varies with a decreasing number of samples available. This is one of the  
166 main conceptual differences compared to the conventional model [18], which does not  
167 consider multicollinearity among tracer properties and uses the mean concentrations  
168 of tracing properties in the available samples. We consider that these assumptions  
169 limit the accuracy of the results and prevent a broader uncertainty analysis.

170 In general, numerical processing computer programs are necessary to calculate the  
171 contributions of each source to target sediments. In this sense, sharing algorithms  
172 and codes that can be tested, modified, and redistributed among researchers and the  
173 general public is of utmost importance to develop new techniques more quickly and  
174 efficiently and improve existing models [2, 17]. Thus, free, open-source, and multi-  
175 platform programming languages that are easy to learn and use are good options for  
176 developing numerical programs that will be available to other researchers, especially  
177 for researchers who have limited programming skills. In this sense, the the growth of  
178 the number of sediment source identification techniques relies on the improvement of  
179 algorithms coded in computer languages that translate mathematical and statistical  
180 modeling into an operational, standardized, and open-source code/language accessible  
181 to the entire scientific community.

182 Despite the challenges raised by studies of the last decades, the sediment fin-  
183 gerprinting technique has proven to be efficient although some further improvement  
184 potential remains possible among the scientific community [15–17, 39]. Many efforts  
185 have been directed at creating protocols to standardize monitoring and modeling tech-  
186 niques, facilitating information and knowledge sharing. Accordingly, the goal of the  
187 current research is to present and make available a computational structure (mod-  
188 ule) to solve over-determined systems applied to sediment tracing (or similar areas)  
189 and systematize the uncertainty analysis and optimization of the number of samples,  
190 which remained under-investigated in the literature during the last years.

191 In this context, the models proposed by Walling and Woodward [18] and Clarke  
192 and Minella [38] will be used to calculate the contributions of each potential source  
193 to suspended sediments collected in rivers. The model outputs will be validated using  
194 a dataset available from the experimental Arvorezinha catchment, in southern Brazil.  
195 The alternative approach proposed in the current research makes it possible to 1) cal-  
196 culate the relative contribution of each source to the composition of the suspended  
197 sediment, allowing us to evaluate the effects of reducing the number of samples and  
198 the associated uncertainties; 2) consider the possible correlations that inherently exist  
199 among the different geochemical variables that compose the set of tracers as deter-  
200 mined from the analysis/use of the variance-covariance matrix in the GLS method;  
201 and 3) calculate the uncertainty variations associated with a change in the number  
202 of samples defining the confidence region of the feasible solutions using the Maha-  
203 lanobis distance. The Python programming language was chosen to implement the  
204 algorithm because it is an open computational language with numerous modules and

1  
2  
3  
4  
5  
6  
7  
8  
9  
205 libraries available for numerical computation, statistics, data processing, and visual-  
206 ization. Libraries such as Pandas, Numpy, Scipy, Tensorflow, and many others are  
207 well-established among data scientists and they allow for rapid prototyping and exper-  
208 imentation. To make this tool available to the entire scientific community, a repository  
209 on the GitHub platform was created to make the Python function module available  
210 for data analysis, visualization, and mathematical routines applied to sediment tracing  
211 modeling.

## 212 2 Material and methods

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
213 The theoretical model applied to the “sediment fingerprinting approach” will be pre-  
214 sented in this section. This approach uses tracers analysed in both soil (i.e., potential  
215 sources) and eroded sediments in multivariate statistical analyses and mathematical  
216 models to optimize over-determined systems. Considering that the suspended sedi-  
217 ments transiting in rivers originate from a mixture of eroded soil from different sources,  
218 the method is based on the principle that the suspended sediment conserves some  
219 bio-physico-chemical characteristics of the sources from which it originates, which is  
220 reflected in its composition. These characteristics that remain “conservative” and dis-  
221 criminant are considered potential tracer properties; we can test their efficiency in  
222 distinguishing sediment sources and determining the contribution of each source to a  
223 given sample collected in the catchment outlet. Nevertheless, various assumptions un-  
224 derlying this approach must be verified and they strongly depend on the number of  
225 samples used to characterize a given tracer (and its spatial variability) in the sources  
226 and eroded sediments.

227 The techniques employed in the fingerprinting approach usually follow three dis-  
228 tinct stages of statistical analysis: 1) range test, 2) discrimination and 3) classification.  
229 The first step refers to determining which geochemical characteristics have the poten-  
230 tial to be selected as tracers among the sources; this step is known as discrimination  
231 analysis. The second step consists of classifying the eroded sediment samples into the  
232  $n$ -dimensional space defined by the source tracer properties. In this step, the relative  
233 contribution of each of the sources to the suspended sediment composition is calculated  
234 by solving a system of over-determined linear equations.

235 The least squares estimation method based on the well-known Gauss-Markov the-  
236 ory has played an essential role in estimating the unknown parameters in linear  
237 regression models [21]. The main assumptions of the technique are that the errors  
238 are assumed to be independent and normally-distributed random quantities with zero  
239 mean and common variance  $\sigma^2$ . In this case, its least squares estimator is the best un-  
240 biased linear estimator of any linear combination of the observations. There are many  
241 ways to define the “best” solution, and one choice is to minimize the sum of squares  
242 of the residuals, where the “best” solution means that the least squares estimators  
243 of its parameters have minimum variance. When there is a certain degree of corre-  
244 lation between the residuals in a regression model, the ordinary least squares (OLS)  
245 and weighted least squares (WLS) may be statistically inefficient or even rely on mis-  
246 leading inferences, in which case it should be preferred to apply the GLS technique  
247 to estimate the unknown parameters in a regression model [23–25]. In addition, the

248 least squares minimization technique allows to estimate regression parameters under  
 249 constraints [19, 20].

## 250 2.1 Review of the classical sediment fingerprinting model

251 The model developed by Yu and Oldfield [40] provides a mathematical formulation to  
 252 obtain the relative contribution of each of the sources to the suspended sediment com-  
 253 position. This model is based on the mass balance of the  $m$  tracers in the different  
 254 potential sources  $g$ . Therefore, the concentrations of  $m$  tracers measured in each sus-  
 255 pended sample sediment are written as the linear combination of the concentrations of  
 256  $m$  tracers measured in each of the different potential sources. The system of equations  
 257 resulting from this expression is given by

$$y_i = \sum_{s=1}^g x_{si} P_s \quad (1)$$

258 where  $i = 1, \dots, m$  the number of tracers. The quantity  $x_{si}$  is the mean concentration of  
 259 the  $i$ -th tracer in the  $s$ -th sediment source and is estimated from samples collected  
 260 in each sediment source;  $P_s$  ( $s = 1, 2, \dots, g$ ) the proportions of sediment supplied by  
 261 the  $g$  sources and  $y_i$  is the concentration of the  $i$ -th tracer in the suspended sediment  
 262 ( $i = 1, \dots, m$ ).

263 Additional equations are assumed to ensure the feasibility of the solutions, since  
 264  $P_1 + P_2 + \dots + P_g = 1$ , with all the  $P_s$  commonly defined as non-negative and lower  
 265 than unity, that is  $0 \leq P_s \leq 1$ ,  $s = 1, \dots, g$ . Furthermore, the samples in each potential  
 266 source and in the suspended sediment are assumed to be statistically independent.

### 267 2.1.1 Constrained least-squares optimization

268 Regression analysis is one of the most widely used statistical tools because it provides  
 269 simple methods for establishing a functional relationship among variables [19–21].  
 270 Constrained optimization procedures to estimate the proportions  $P_s$  by least-squares  
 271 objective function (OLS) were first used by Walling and Woodward [18]. The authors  
 estimated the  $P_s$  by minimization of

$$\sum_{i=1}^m \left\{ \left( y_i - \left( \sum_{s=1}^g x_{si} P_s \right) \right) / y_i \right\}^2 \quad (2)$$

Subject to the constraints

$$0 \leq P_s \leq 1 \quad (3)$$

and

$$\sum_{s=1}^g P_s = 1 \quad (4)$$

268 The division by  $y_i$  ( $i = 1, \dots, m$  the number of tracers) provides a type of scaling  
 269 to data for different tracers, which may differ greatly in variability.

270 It should be noted that, provided the number of source areas ( $g$ ) is not greater than  
 271 the number of tracers ( $m$ ), criterion (2) could in principle be minimized by the least

272 squares, even if there was only one sample from each sediment source and only one  
 273 sample of suspended sediment. In this case, an objective function is used to minimize  
 274 Eq. (2).

275 This model has provided the main strategy applied by researchers to estimate the  
 276 relative contribution of each of the sources to the sediment composition. However, the  
 277 approach is associated with some limitations regarding the influence of the number of  
 278 samples used to conduct it, since it considers the mean value of the sample properties  
 279 and it does not take into account collinearity between the variables used. Both limi-  
 280 tations have a strong impact on the quality of the results and prevent an uncertainty  
 281 analysis.

282 Basead on the model proposed by Walling and Woodward [18], Clarke [36] and  
 283 Clarke and Minella [38] applied the OLS to the original data (OLS\_Clarke model)  
 284 and introduced the constraints on solving the system, following the steps below:

1º) The division of the terms of the Equation 1 by  $y_i$ :

$$1 = \sum_{s=1}^g \left( \frac{x_{si}}{y_i} \right) P_s = \left( \frac{x_{1i}}{y_i} \right) P_1 + \left( \frac{x_{2i}}{y_i} \right) P_2 + \dots + \left( \frac{x_{gi}}{y_i} \right) P_g$$

So, the equation is rewritten:

$$1 = \sum_{s=1}^g a_{si} P_s \quad (5)$$

285 where  $a_{si} = \frac{x_{si}}{y_i}$ .

Nominating  $\mathbf{W}$  the matrix generated by  $\mathbf{a}_s$ ,  $s = 1 \dots g$  the column vectors ( $i =$   
 $1, \dots, m$  the number of tracers),  $\mathbf{p}$  column vectors of  $P_s$  and  $\mathbf{1}$  column vectors of 1,  
 then the underdetermined system of linear equations (Eq. 5) is written matricial form

$$\mathbf{W}\mathbf{p} = \mathbf{1} \quad (6)$$

286 2º) An ordinary least squares method (OLSE) is equivalent to solving the minimiza-  
 287 tion problem

$$\min_p (\mathbf{W}\mathbf{p} - \mathbf{1})^T (\mathbf{W}\mathbf{p} - \mathbf{1}), \quad (7)$$

288 The associated normal system to Eq. 6 is given by

$$\mathbf{W}^T \mathbf{W}\mathbf{p} = \mathbf{W}^T \mathbf{1}, \quad (8)$$

289 3º) Adding the constraints  $P_1 + P_2 + \dots + P_g = 1$  to the normal system of equations  
 290 associated with Eq. (8) results as Eq. (9):

$$\mathbf{A}\mathbf{P} = \mathbf{Z} \quad (9)$$

where  $\mathbf{A}$  is  $(g+1) \times (g+1)$  matrix,  $\mathbf{P}$  is  $(g+1) \times 1$  vector and  $\mathbf{Z}$  is  $(g+1) \times 1$  vector are given, respectively, by

$$\mathbf{A} = \begin{pmatrix} \langle \mathbf{a}_1, \mathbf{a}_1 \rangle & \langle \mathbf{a}_1, \mathbf{a}_2 \rangle & \cdots & \langle \mathbf{a}_1, \mathbf{a}_g \rangle & 1 \\ \langle \mathbf{a}_2, \mathbf{a}_1 \rangle & \langle \mathbf{a}_2, \mathbf{a}_2 \rangle & \cdots & \langle \mathbf{a}_2, \mathbf{a}_g \rangle & 1 \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ \langle \mathbf{a}_g, \mathbf{a}_1 \rangle & \langle \mathbf{a}_g, \mathbf{a}_2 \rangle & \cdots & \langle \mathbf{a}_g, \mathbf{a}_g \rangle & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{pmatrix} \quad (10)$$

$$\mathbf{P} = \begin{pmatrix} P_1 \\ P_2 \\ \vdots \\ P_g \\ 1 \end{pmatrix} \quad (11)$$

$$\mathbf{Z} = \begin{pmatrix} \langle \mathbf{a}_1, \mathbf{1} \rangle \\ \langle \mathbf{a}_2, \mathbf{1} \rangle \\ \vdots \\ \langle \mathbf{a}_g, \mathbf{1} \rangle \\ 1 \end{pmatrix} \quad (12)$$

where  $\mathbf{a}_s$  for  $s = 1 \cdots g$  are the dimensionless column vectors of  $\mathbf{W}$  matrix.  
 4<sup>o</sup>) Then, Equation 9 must be solved.

$$\hat{\mathbf{P}} = \mathbf{A}^{-1} \mathbf{Z} \quad (13)$$

A solution of Equation 9 will automatically satisfy the constraint  $P_1 + P_2 + \cdots + P_g = 1$ , although it may not satisfy the inequality constraints  $0 < P_1, P_2, \cdots, P_g < 1$ . In this case, the solution must be tested in order to satisfy the inequality to obtain all the feasible solutions.

In this approach, the authors were able to analyze the influence of reducing the number of samples used in the uncertainty analysis. However, the model did not take into account the collinearity between the potential tracing variables.

### 2.1.2 The GLS\_Clarke model

The method developed by Clarke and Minella [38] presents a way to calculate the uncertainty of the sources apportionment when the number of samples of sources and/or target varies. Besides, the multicollinearity between potential tracing properties is incorporated from the variance-covariance matrix (GLS). For the GLS\_Clarke model, the authors suggest the following procedure:

1<sup>o</sup>) The dimensional underdetermined system

$$y_i = \sum_{s=1}^g x_{si} P_s \quad (14)$$

is written in matrix by naming  $\mathbf{X}$  as the matrix generated by  $\mathbf{x}_s$ ,  $s = 1 \cdots g$  the column vectors ( $m$ -dimensional vectors of the tracers),  $\mathbf{p}$  column vectors of  $P_s$  and  $\mathbf{Y}$  column vectors of concentration in suspended sediment. Then, the under-determined system of linear equations (Eq. 14) is given by

$$\mathbf{X}\mathbf{p} = \mathbf{Y} \quad (15)$$

306 **2°)** A generated least squares method (GLSE) is equivalent to solving the minimiza-  
307 tion problem

$$\min_p (\mathbf{X}\mathbf{p} - \mathbf{Y})^T \mathbf{S}^{-1} (\mathbf{X}\mathbf{p} - \mathbf{Y}), \quad (16)$$

308 where  $\mathbf{S}$  is the ( $m \times m$ ) variance-covariance matrix of the "dependent" variable  $y_i$ ,  
309  $i = 1 \cdots m$ .

310 The associated normal system to Eq. 15, which is known as the Aitken equations  
311 ([22]) is given by

$$\mathbf{X}^T \mathbf{S}^{-1} \mathbf{X} \mathbf{p} = \mathbf{X}^T \mathbf{S}^{-1} \mathbf{Y}, \quad (17)$$

312 **3°)** Adding the constraints  $P_1 + P_2 + \dots + P_g = 1$  to the normal system of equations  
313 associated Eq. (17) results into:

$$\mathbf{B}\mathbf{p} = \mathbf{V} \quad (18)$$

where  $\mathbf{B}$  is  $(g+1) \times (g+1)$  matrix,  $\mathbf{p}$  is  $(g+1) \times 1$  vector (Eq. 11) and  $\mathbf{V}$  is  $(g+1) \times 1$  vector are given, respectively, by

$$\mathbf{B} = \begin{pmatrix} \langle \mathbf{x}_1, \mathbf{b}_1 \rangle & \langle \mathbf{x}_1, \mathbf{b}_2 \rangle & \cdots & \langle \mathbf{x}_1, \mathbf{b}_g \rangle & 1 \\ \langle \mathbf{x}_2, \mathbf{b}_1 \rangle & \langle \mathbf{x}_2, \mathbf{b}_2 \rangle & \cdots & \langle \mathbf{x}_2, \mathbf{b}_g \rangle & 1 \\ \vdots & \ddots & & \vdots & \\ \langle \mathbf{x}_g, \mathbf{b}_1 \rangle & \langle \mathbf{x}_g, \mathbf{b}_2 \rangle & \cdots & \langle \mathbf{x}_g, \mathbf{b}_g \rangle & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{pmatrix} \quad (19)$$

$$\mathbf{V} = \begin{pmatrix} \langle \mathbf{x}_1, \mathbf{y} \rangle \\ \langle \mathbf{x}_2, \mathbf{y} \rangle \\ \vdots \\ \langle \mathbf{x}_g, \mathbf{y} \rangle \\ 1 \end{pmatrix} \quad (20)$$

where  $\mathbf{b}_s$  and  $\mathbf{y}$  for  $s = 1 \cdots g$  are the column vectors of the product of matrices  $\mathbf{S}^{-1}$  by  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively. In that way, entries of vectors are given by

$$b_{jk} = \sum_{l=1}^m s_{jl} \alpha_{lk} \quad (21)$$

314 and

$$y_j = \sum_{l=1}^m s_{jl} y_l \quad (22)$$

for  $j = 1 \dots m$  and  $k = 1 \dots g$ .  
 4<sup>o</sup>) Equation 18 can then be solved:

$$\hat{\mathbf{P}} = \mathbf{B}^{-1} \mathbf{V} \quad (23)$$

where  $\hat{\mathbf{P}}$  is the vector solution to the sediment source contributions.

### 2.1.3 Mahalanobis distance

The Mahalanobis distance is a multivariate distance metric that measures the distance between a point (vector) and a distribution. This distance differs from the Euclidean distance because it is calculated using the inverse of the variance-covariance matrix of the dataset [41, 42]. It is a useful metric since it measures the distances taking into account the correlation between the variables, even if the data are not normalized values (different scale). The Mahalanobis distance was used to determine the confidence region and estimate the associated uncertainties. Let  $n^*$  be the number of feasible solutions of the systems (13 or 23), thereby for each  $j = 1 \dots n^*$  there is the solution vector  $\mathbf{P}_j = [P_1 P_2 \dots P_{g-1}]^T$ , and a vector of the averages of the feasible solutions  $\overline{\mathbf{P}}^* = [\overline{P_1} \overline{P_2} \dots \overline{P_{g-1}}]^T$ , both of dimension  $(g-1) \times 1$ . Thus, the Mahalanobis distance is defined by:

$$d_j^2 = (\mathbf{P}_j - \overline{\mathbf{P}}^*)^T \Sigma^{-1} (\mathbf{P}_j - \overline{\mathbf{P}}^*) \quad (24)$$

where  $\Sigma$  is the variance-covariance matrix  $(g-1) \times (g-1)$  of the feasible solutions  $n^*$ .

## 2.2 Case study: Arvorezinha experimental catchment

The dataset used for the evaluations of the mathematical models (OLS\_Clarke model and GLS\_Clarke model) explored in the Sections Sections 2.1.1 and 2.1.2 was taken from a sediment tracing study performed in the Arvorezinha experimental catchment between 2002 and 2006 (Figure 4). This catchment is located in southern Brazil (28°52' S and 52°05' O), it covers a surface area of 1.19 km<sup>2</sup>, and is located on the edge of the Brazilian southern plateau. In its upper part, the topography is gently rolling, and in its lower parts, it is characterized by shorter and steeper slopes. Volcanic rocks and shallow and fragile soils (Entisols and Inceptisols) characterize the geology and the pedology of the catchment. The climate is subtropical super-humid meso-thermic (i.e., Cfb). The mean annual precipitation is 1605 mm over 50-yr period, evenly distributed throughout the year. Land use is mainly agricultural, with much of the land used for tobacco cultivation. Soil erosion is the main soil degradation process thereby generating high sediment yields Minella et al. [43].

Insert Figure 1 here

This experimental catchment is used to investigate hydrological and erosion processes at the catchment scale in this region of South America. The primary goal of

350 this research was monitoring liquid and solid discharges during significant rainfall-flow  
351 events. This monitoring data is then used to improve water and sediment transfer mod-  
352 eling techniques, and to identify and quantify sediment source contributions through  
353 the implementation of sediment fingerprinting techniques [3, 44]. In the study of  
354 Minella et al. [43], three potential sediment sources were considered, and their relative  
355 contributions to 24 suspended sediment samples collected during 20 significant rainfall-  
356 flow events were quantified. The sediment sources evaluated were channel banks (*CB*),  
357 unpaved roads (*UR*) and crop fields (*CF*). For tracer determination, 9 *CB* samples, 9  
358 *UR* samples and 20 *CF* soil samples were collected across the catchment. The set of  
359 62 samples of suspended sediments and potential sources were characterized for their  
360 elemental geochemistry (total concentrations in P, K, Ca, Na, Mg, Cu, Pb, Cr, Co, Zn,  
361 Ni, Fe, Mn, and Al). As a first step, the individual analysis of the discriminant ability  
362 of each element was performed by means of a range test to evaluate the conservativity  
363 and the Kruskal-Wallis test to estimate the discrimination power [7]. Subsequently,  
364 the best set of tracer elements was determined based on multivariate or discriminant  
365 functional analyses to define the best set of variables for discriminating sources. This  
366 analysis maximizes the discrimination between the sources and minimizes the number  
367 of variables required. The method is based on minimizing the Wilks' Lambda index  
368 ( $\lambda^*$ ), a component of the multivariate analysis of variance.

369 The methods of tracer selection are fundamental in sediment fingerprinting because  
370 it may decrease the deleterious effect of multicollinearity. Davis and Sampson [45]  
371 reported that redundant variables may weaken the analysis due to the reduction of the  
372 degrees of freedom of the errors and may affect the feasibility the inversion operation  
373 of the variance-covariance matrix. In addition, they reduce the dimensionality of the  
374 problem, significantly improving the ability of the model to find a solution and the  
375 associated errors. Seven ( $m = 7$ ) of the 14 chemical elements previously analyzed were  
376 selected to compose the optimal data set in this case study (Fe, Mn, Cu, Zn, Ca, K  
377 and P are the optimal tracers).

### 378 **2.3 Fingerprinting computational framework applied to** 379 **Arvorezinha catchment data set**

380 The Python programming language was used to implement the OLS\_Clarke and  
381 GLS\_Clarke models. The samples collected from  $g$  sources and suspended sediments,  
382 with the respective concentrations of  $m$  tracers with  $g \leq m$ , provide the input data of  
383 the algorithm. It consists of a module and external functions used repeatedly for the  
384 approximate resolution of the over-determined linear system with  $m + 1$  equations and  
385  $g$  unknowns in the form  $y_i = \sum x_{si}P_s$  with  $s = 1, \dots, g$  and  $i = 1, \dots, m$  conditioned on  
386  $1 = P_1 + P_2 + \dots + P_g$ . From the solutions obtained, the feasible proportions provide a  
387 means to estimate the uncertainty of the results as a function of reducing the number  
388 of samples used. The process is repeated 100 times, generating a cloud of feasible  
389 solutions in the  $P_1P_2$  plane. A flowchart was developed to provide an overview of the  
390 structural organization of the algorithm.

391  
392 Insert Figure 2 - Flowchart here  
393

Each block in the flowchart has tasks executed by commands and functions described below:

1. Python import library: The import statement Python.
2. Data set: The data representing each  $g$  source and suspended sediment should be organized in  $g + 1$  spreadsheets, where each column represents one of the optimal tracers and each row represents a particular sample. For instance, spreadsheet  $CB$  will be a  $9 \times 7$  table that corresponds to 9 samples collected in channel banks and 7 optimal tracers analysed in these, spreadsheet  $UR$  will then be a  $9 \times 7$  table that contains 9 samples collected unpaved roads and 7 optimal tracers, spreadsheet  $CF$  will finally correspond to a  $20 \times 7$  table that contains 20 samples collected crop fields and 7 optimal tracers, and spreadsheet  $Y$  will include a  $24 \times 7$  table that associates 24 suspended sediment samples and their 7 optimal tracers. These subsets are saved in four matrices that will then be used in the next operations.
3. Subset random: It randomly chooses, without repetition, the samples from each subset; in this case, we can select  $1 \leq n_{CB} \leq 9$ ,  $1 \leq n_{UR} \leq 9$ ,  $1 \leq n_{CF} \leq 20$  and  $1 \leq n_Y \leq 24$ . Hence, the total number of possible combinations involving the samples from each set for the coefficients of Equation (1), is equal to  $C_n = n_Y \times n_{CB} \times n_{UR} \times n_{CF}$ , which will be reduced as the number of samples used for analysis decreases. For example, we can choose the sequence  $\{Y_n\}_{n=1}^6 = \{24, 20, 16, 12, 8, 4, 2\}$  as the number of suspended sediment samples  $Y$  and keep the number of the other subsets fixed, then the sequence of possible combinations is given by  $\{C_n\}_{n=1}^6 = \{38880, 25920, 19440, 12960, 6480, 3240\}$ . Each selected  $CB$ ,  $UR$  and  $CF$  sample will correspond to a column of the source matrix ( $W$  or  $X$ ) and  $Y$  to a column vector of the suspended sediment.
4. Models: For each of the drawings, the external functions `OLS_Clarke` and `GLS_Clarke` described in Sections 2.1.1 and 2.1.2 are called.
5. Cloud Proportions: It calls a third external function to compute the proportions  $P_1, P_2, \dots, P_g$  of the  $g$  sources. Among the solutions obtained in the `OLS_Clarke` and `GLS_Clarke` procedures, the algorithm retains those that are feasible (i.e., such that  $0 < P_1, P_2, P_3 < 1$ ). The process is then repeated 100 times, generating a cloud of feasible solutions in the  $P_1P_2$  plane.
6. Confidence regions: The confidence interval region (95%) is calculated from the set of feasible solutions  $n^*$  for each  $n_Y$  or  $n_{CF}$  using the procedure described in Section 2.1.3. Lastly, a function generates the graphical visualization of this region.
7. Outputs: The calculation of the mean of the areas of the confidence region, standard deviation, number of combinations, number of feasible solutions, mean of the values of  $P_1, P_2, P_3$ , and the coefficient of variation.

The Fingerprinting repository on the GitHub platform was created to provide a module of functions in Python for data analysis, visualization, and mathematical routines applied to sediment tracing modeling. The model proposed by [38] is implemented and can be modified and applied to different databases. The module consists of functions in Python and has the libraries Numpy [46], Scipy [47], and Matplotlib [48] as the dependencies. A Jupyter Notebook file provides examples of using some functions and reproducing the results published in this work. From

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

439 this example, we can explore other data sets, modify input parameters in the func-  
440 tions, operationalize models, and create graphs. The repository can be accessed at  
441 <https://github.com/tiagoburiol/Fingerprinting>.

### 442 3 Results and discussion

443 With the algorithm implemented in Python, it is possible to execute the sequence of  
444 instructions that operationalize the OLS\_Clarke and GLS\_Clarke models and the  
445 uncertainty calculation. The simulations allow the user to compare the results gener-  
446 ated by both methods to solve over-determined systems, select solutions utilizing the  
447 constraints, calculate the contributions of each source to the sediments, and statistical  
448 parameters such as the confidence interval region, the standard deviation, and the co-  
449 efficient of variation. The user can also indicate the number of samples in the different  
450 groups of samples (target sediments or soil sources) to analyze the effect of reducing  
451 their number on the calculations of the  $P_s$  values and the impact of this reduction on  
452 the uncertainty associated with the output results.

453 From the simulation results, we quantified the influence of collinearity on the results  
454 of the source ratios. In the tables and figures below, the results for the OLS\_Clarke  
455 model correspond to the application of the model without considering the effect of  
456 collinearity, and the results for the GLS\_Clarke model, which indicate the application  
457 of the model when considering the effect of collinearity. The comparison between the  
458 two models is obtained simultaneously, that is, using the same subset of randomly  
459 drawn samples.

460 In addition to the effect of collinearity, we also present the analysis of the in-  
461 creased uncertainties associated with a reduction in the number of samples from  
462 both the suspended sediment samples ( $nY$ ) and the crop fields samples ( $nCF$ ). The  
463 source  $CF$  was choose the higher number of data. The two simulation to evaluate the  
464 sample number reduction were performed using the OLS\_Clarke and GLS\_Clarke  
465 models under constraints described in Sections 2.1.1 and 2.1.2. The values of  $P_1$ ,  
466  $P_2$ , and  $P_3$ , which correspond to the contributions from the  $CB$ ,  $UR$ , and  $CF$   
467 sources, are identified by: a) Simulation 1: reduction of suspended sediment samples  
468  $nY = \{24, 20, 16, 12, 8, 4, 2\}$  (Table 1) and b) Simulation 2: reduction of the number  
469 of crop samples  $nCF = \{20, 16, 12, 8, 4, 2\}$  (Table 2).

470  
471 Insert TABLE 1 here

472  
473 Insert TABLE 2 here

474  
475 The first result is that the inclusion of the variance-covariance matrix, and conse-  
476 quently the effect of collinearity, affects the results of both simulations. The significant  
477 difference in the control parameters (mean area,  $\sigma$  and  $n^*$ ) and in the mean values of  
478  $P_1$ ,  $P_2$ , and  $P_3$  can also be observed.

479 Although both models indicate the  $CF$  source ( $\overline{P}_3$ ) as the main source of sediment,  
480 the OLS\_Clarke model suggests that the contribution from the  $CB$  source ( $\overline{P}_1$ ) is  
481 greater than that from the  $UR$  source ( $\overline{P}_2$ ). Nonetheless, the GLS\_Clarke model

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

482 indicates that the *UR* source ( $\overline{P_2}$ ) supplies more sediment than the *CB* source ( $\overline{P_1}$ ),  
483 corroborating field evidence and observations. The catchment is characterised by a  
484 dense road network, where many sections are inadequately allocated given the local  
485 topography and without the implementation of runoff control practices. According to  
486 a study by [49] in this catchment, unpaved roads are severely affected by water erosion  
487 and contribute significantly to sediment yield.

488 These results suggest that it is important to take into consideration the possible  
489 correlations between the chemical variables that comprise the tracer set as the GLS  
490 method performed. In this model, the values of  $n^*$  are lower, and the value of  $\sigma$   
491 decreases significantly, although the average area of the confidence region is larger.

492 The comparison between both methods demonstrates the interest of extracting the  
493 effect of collinearity existing in the set of tracers even if these have been appropriately  
494 selected in the discrimination phase by the range test, KruskalWallis test, and min-  
495 imization of Wilk's lambda. According to Johnson et al. [50] the variance-covariance  
496 matrix allows to extract the effects of correlation among the variables in the solution  
497 of the systems as those presented in this work. De Maesschalck et al. [41] through the  
498 Mahalanobis distance, also demonstrated the applicability of this approach in reducing  
499 the uncertainties associated with collinearity.

500 Figure 3 shows the point distributions that express the feasible solutions in the  
501 95% confidence interval region. In the GLS\_Clarke model, the point distributions  
502 show a higher "density" in the region near the mean value than those obtained from  
503 the OLS\_Clarke model.

504  
505 Insert Figure 3 here  
506

507 However, it is important to note that the possibility of considering collinearity  
508 among the tracer variables does not diminish the importance of previous discrimi-  
509 nation analyses in defining the best set of tracers. Including the variance-covariance  
510 matrix of the previously selected variables considers the existing correlation between  
511 the chosen variables, even if it remains limited. It is known that sediment production  
512 in catchments can result from different processes (diffuse or concentrated, agricul-  
513 tural or fluvial, superficial or deep erosion) in which different variables will reflect the  
514 set of operating processes. Therefore, we consider it fundamental to analyse a signifi-  
515 cant amount of tracers that maximize the discriminating capacity, which will offer the  
516 physico-chemical basis to differentiate the sources.

517 Even after the selection analysis of the set of variables in the minimization of  
518 the Wilks' Lambda index ( $\lambda^*$ ), a certain degree of collinearity may influence the  
519 final result. Moreover, to exclude the variables (tracers) that present some degree of  
520 collinearity would be neglecting part of the useful existing variability. This fact can be  
521 verified in the higher uncertainties by the simulations of the OLS\_Clarke model that  
522 does not consider the variance-covariance matrix in the solution of the over-determined  
523 equation systems.

524 In Simulation 1 and for both models, the number of feasible solutions and the  
525 mean area of the confidence region decreased as the number of suspended sediment  
526 samples considered decreased, corroborating the result of Clarke and Minella [38].

527 The distribution of points expressing the feasible solutions within the confidence  
528 interval region (95%) for the cases  $nY = 12$  and  $nY = 4$  are illustrated in Figures 4  
529 and 5, respectively. We can observe the lower number of feasible solutions. Besides  
530 this, and, as in Figure 3, the points are more concentrated around the  $P_1$  and  $P_2$   
531 mean values in the GLS\_Clarke model.

532  
533 Insert FIGURE 4 here

534  
535 Insert FIGURE 5 here

536  
537 Figure 6 shows the coefficients of variation ( $CV$ ) obtained with the reduced number of  
538 sediment samples for each model (OLS\_Clarke and GLS\_Clarke). In these  
539 plots, we can observe the impact of reducing the number of eroded sediment samples  
540 collected in the river when defining the  $P_s$  values. For all sample quantities, the  
541 standard deviation is lower with the GLS\_Clarke model. Additionally, the difference  
542 in the  $CV$  rates of change in the OLS\_Clarke model is much higher in a few samples  
543 (e.g., in the range below ten samples).

544  
545 Insert FIGURE 6 here

546  
547 Figures 7, Fig. 8, and Fig.9 show the regions of the confidence interval (95%)  
548 with the reduction in the number of samples from  $nCF = 20$ ,  $nCF = 12$ , and  
549  $nCF = 4$ , respectively. Notably, the variance-covariance matrix in the GLS\_Clarke  
550 model comes from the  $CF$  source data. Therefore, the point distributions in Figures  
551 6 (b) and 7 (b) differ even though both simulations use the total number of existing  
552 samples (ie  $nY = 24$  and  $nCF = 20$ ).

553  
554 Insert FIGURE 7 here

555  
556 Insert FIGURE 8 here

557  
558 Insert FIGURE 9 here

559  
560 As for the previous scenario simulation, the number of feasible solutions decreases  
561 as the number of suspended sediment samples considered decreases, as does the  
562 average area of the confidence region. The cases simulated by both models have  
563 mean proportions of  $\overline{P_1}$  greater than the mean proportions of  $\overline{P_2}$ , indicating that the  
564 relative contribution to the  $CB$  source is greater than that of  $UR$  in the suspended  
565 sediment composition ( $\overline{P_1} > \overline{P_2}$ ). This result differs from that of Simulation 1, in  
566 which the GLS\_Clarke model showed similar solutions as the OLS\_Clarke model  
567 with  $\overline{P_2} < \overline{P_1}$ . This indicates that the GLS method responds directly to the choice  
568 of inclusion of the variance-covariance matrix. Figure 10 presents the  $CV$  values  
569 associated with the reduction in the number of crop fields samples considered, and  
570 Table 2 lists the mean values of the proportions obtained in such a situation.

572 Insert FIGURE 10 here

573  
574 The OLS\_Clarke model showed higher  $CV$  values than the GLS\_Clarke model  
575 when analysing the uncertainties associated with the results obtained with a reduced  
576 number of samples, indicating, once again, that the uncertainties decrease as the effect  
577 of collinearity among the variables is considered in the system over-determined by the  
578 presence of the variance-covariance matrix.

579 In all simulated cases, the relative contribution of  $\overline{P}_3$  is greater, indicating that the  
580 crop fields ( $CF$ ) source provides the main sediment source in the investigated catch-  
581 ment. This result corroborates those of other previous studies conducted by Minella  
582 et al. [44] in the catchment and all documented field observations.

## 583 4 Conclusions

584 This study presented an updated and open access computational framework for solving  
585 over-determined systems applied to sediment tracing based on the models previously  
586 proposed by Clarke and Minella [38]. The algorithm enabled us to analyze and com-  
587 pare the use of different metrics, optimize the procedure for calculating the area that  
588 expresses the degree of uncertainty associated with the number of samples taken into  
589 account to document source and target sediment properties, and consider the possible  
590 correlations that naturally exist between the different variables that compose the set  
591 of tracers used for sediment tracing.

592 Our findings confirmed the assumption of increased uncertainty as the number of  
593 samples considered decreases either for the potential sources or the target sediment  
594 samples. Moreover, considering the variance-covariance matrix to find the solution of  
595 the over-determined system allowed us to consider the deleterious effects of collinearity  
596 between tracers in sediment tracing studies. The implemented algorithm allowed us  
597 to compare the two modeling strategies and simulate multiple scenarios of sample  
598 number reduction.

599 The Python language enables the use of an easy-to-manipulate code, and will facili-  
600 tate the shared implementation of the model and the inclusions of model modifications  
601 that may be suggested by the research community to keep improving and standardiz-  
602 ing the sediment tracing protocol. With this tool, new perspectives are opened to to  
603 provide a helping decision tool to define the number of samples required to charac-  
604 terise the potential sources and/or sediment based on the uncertainty analysis of the  
605 set of samples available, which is fundamental for the advancement of research in en-  
606 vironmental monitoring and modeling, as well as for the effective management of soil  
607 and water resources at the catchment scale.

608 Considering that methodological development occurs through the contribution of  
609 professionals from different areas of knowledge, exchanging ideas and the efficient use  
610 of computing resources promotes the standardization and accuracy of techniques. In  
611 this context, the sharing of the algorithm implemented in Python provides speed and  
612 transparency in scientific development and increases access to new methods related to  
613 fingerprinting, providing researchers with an algorithm that enables them to evaluate

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

614 uncertainties in reducing the number of samples and the influence of the collinearity  
615 of the set of tracers, which can also be adapted to other areas of knowledge.

## 616 **Declarations**

### 617 **Ethical Approval**

618 Not applicable

### 619 **Consent to Participate**

620 Not applicable

### 621 **Consent to Publish**

622 The authors hereby consents to publication of the manuscript in the Computational  
623 and Applied Mathematics journal.

### 624 **Authors Contributions**

625 Conceptualization: Lidiane Buligon and Jean Paolo Gomes Minella; Methodology:  
626 Jean Paolo Gomes Minella, Tiago Martinuzzi Buriol, and Lidiane Buligon; Formal  
627 analysis and investigation: Jean Paolo Gomes Minella, Tiago Martinuzzi Buriol, and  
628 Lidiane Buligon; Writing - original draft preparation: Lidiane Buligon, Jean Paolo  
629 Gomes Minella, and Tiago Martinuzzi Buriol; Writing - review and editing: Olivier  
630 Evrard; Funding acquisition: Lidiane Buligon, Jean Paolo Gomes Minella, and Olivier  
631 Evrard.

### 632 **Funding**

633 CNPq (National Council for Scientific and Technological Development) Process  
634 200008/2023-4. CAPES (Coordination for the Improvement of Higher Education  
635 Personnel) Process 88887.310201/2018-00. CEA (Atomic Energy Commission).

### 636 **Competing Interests**

637 The authors declare that they have no conflict of interest that are relevant to the  
638 content of this article.

### 639 **Availability of data and materials**

640 Data will be made available on reasonable request.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

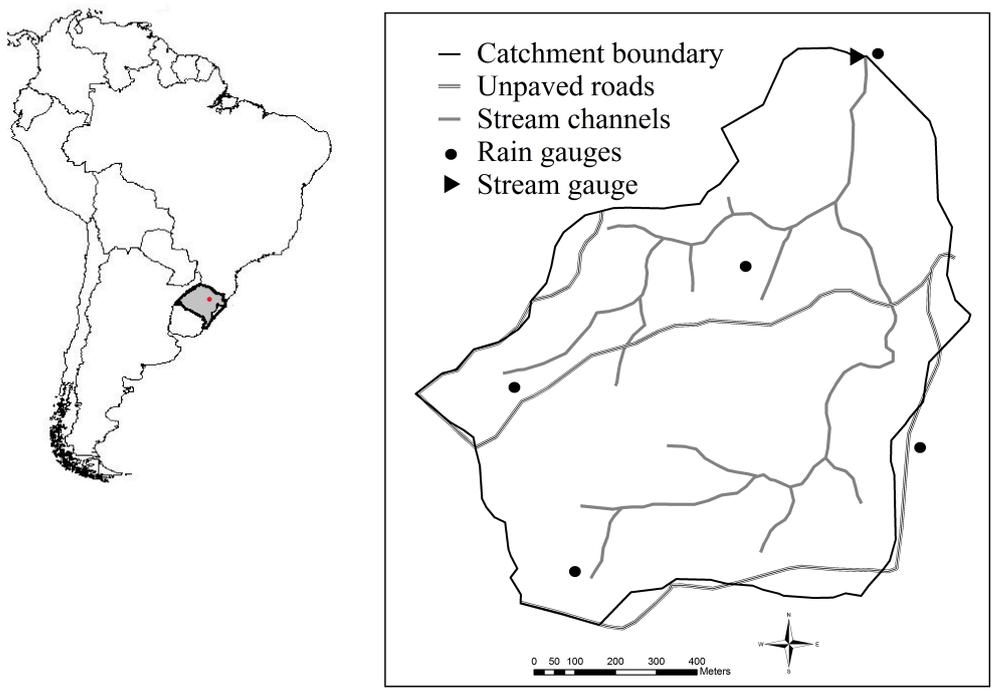


Figure 1 Arvorezinha experimental catchment: Adapted from [3].

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

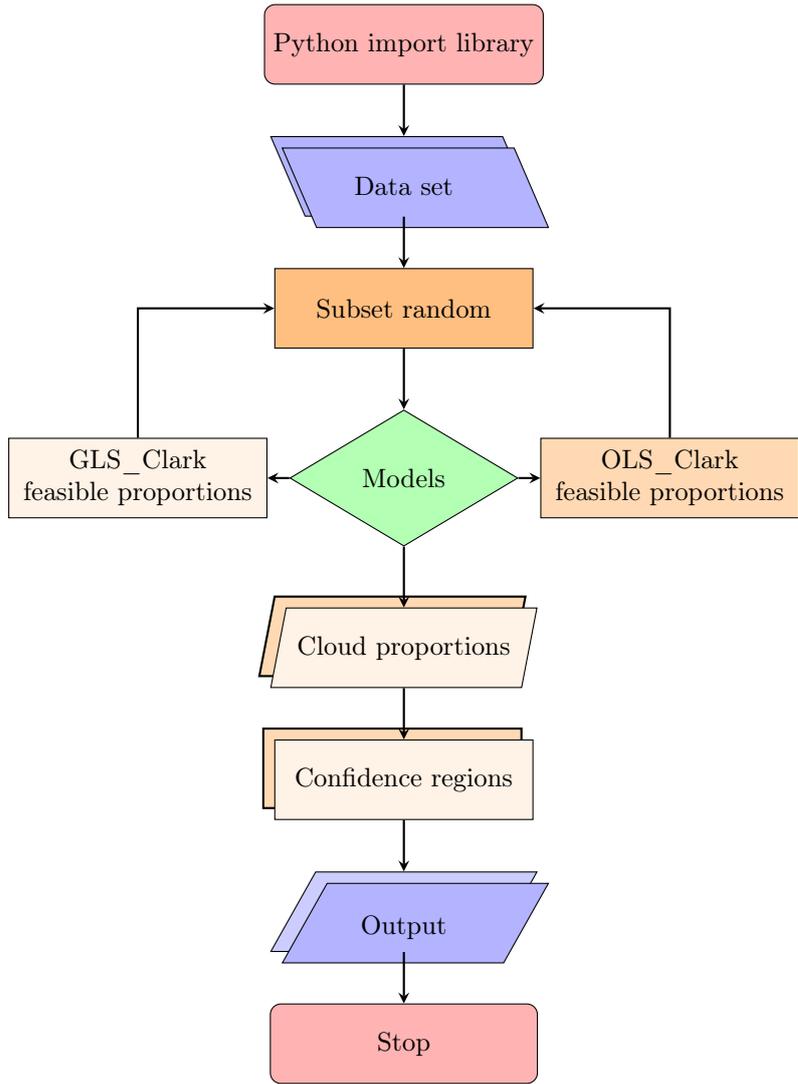
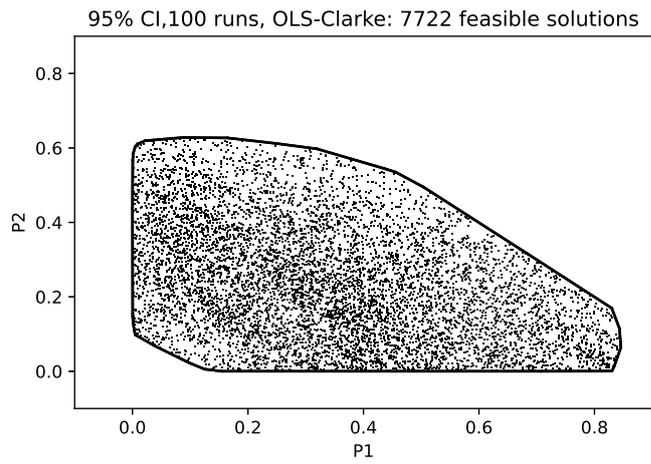
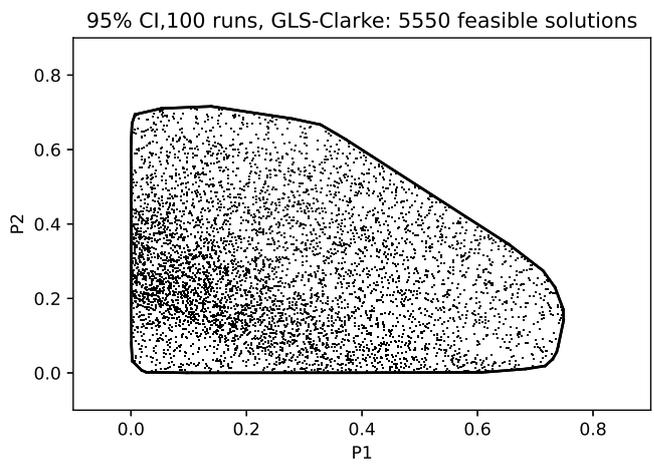


Figure 2 Flowchart of the structural organization of the algorithm.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



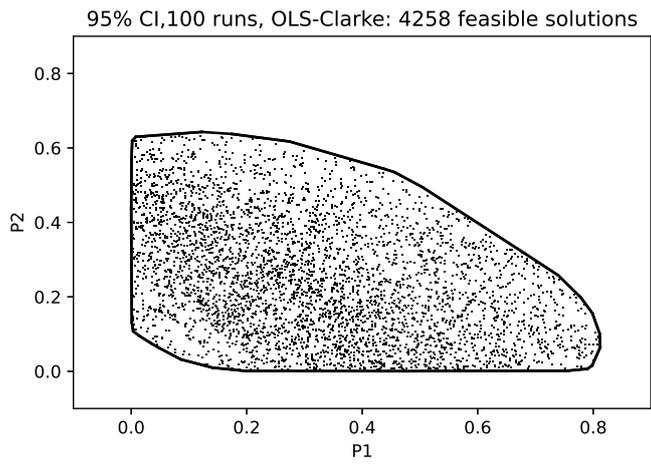
(a) OLS\_Clarke



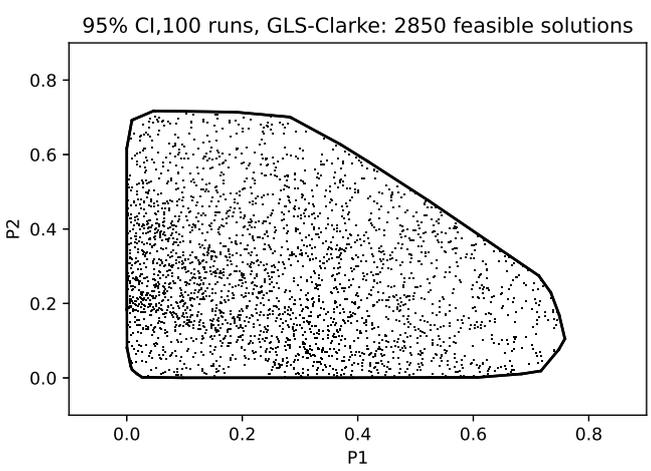
(b) GLS\_Clarke

**Figure 3** The feasible solutions in the 95% confidence region. The mean of proportions  $(\bar{P}_1, \bar{P}_2) = (0.341, 0.238)$  and  $(\bar{P}_1, \bar{P}_2) = (0.259, 0.276)$  obtained by a) OLS\_Clarke model and b) GLS\_Clarke model, respectively. Considering  $nY = 24$  over 100 repetitions.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



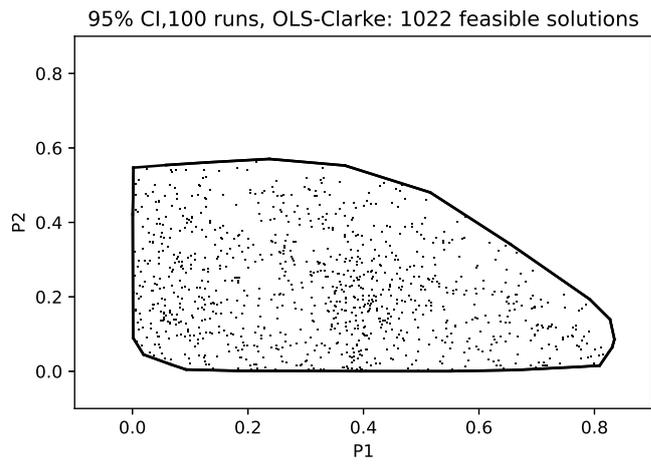
(a) OLS\_Clarke



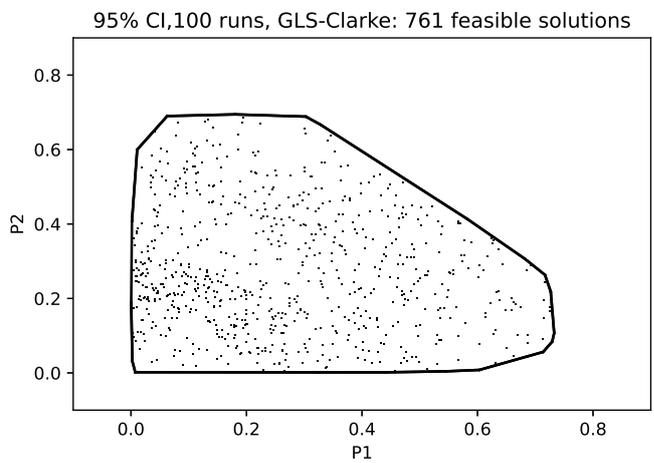
(b) GLS\_Clarke

**Figure 4** The feasible solutions in the 95% confidence region. The mean of proportions  $(\bar{P}_1, \bar{P}_2) = (0.319, 0.248)$  and  $(\bar{P}_1, \bar{P}_2) = (0.267, 0.278)$  obtained by a) OLS\_Clarke model and b) GLS\_Clarke model, respectively. Considering  $nY = 12$  over 100 repetitions.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



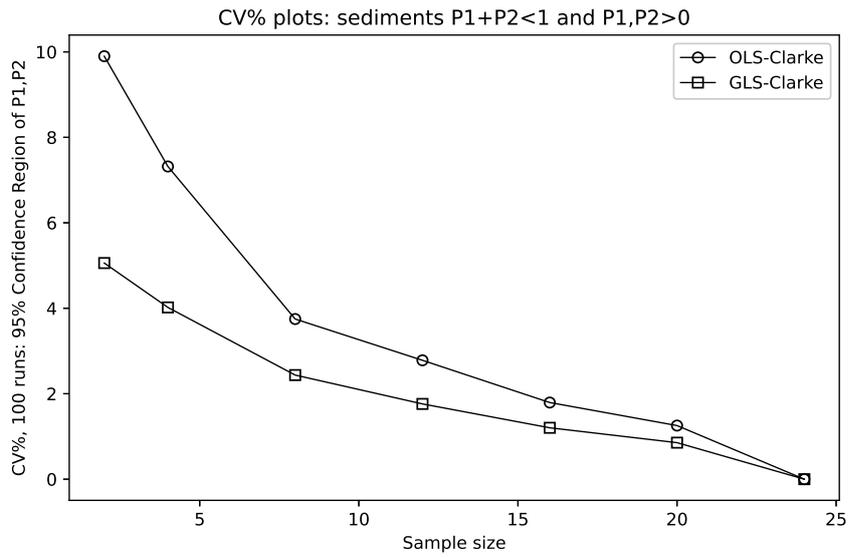
(a) OLS\_Clarke



(b) GLS\_Clarke

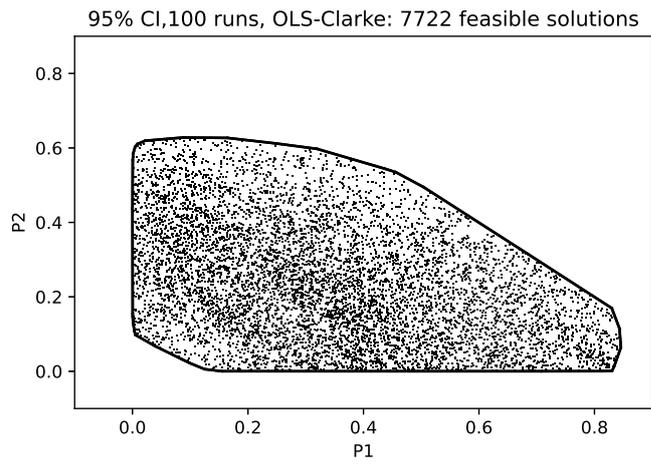
**Figure 5** The feasible solutions in the 95% confidence region. The mean of proportions  $(\bar{P}_1, \bar{P}_2) = (0.333, 0.219)$  and  $(\bar{P}_1, \bar{P}_2) = (0.266, 0.276)$  obtained by a) OLS\_Clarke model and b) GLS\_Clarke model, respectively. Considering  $nY = 4$  over 100 repetitions.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

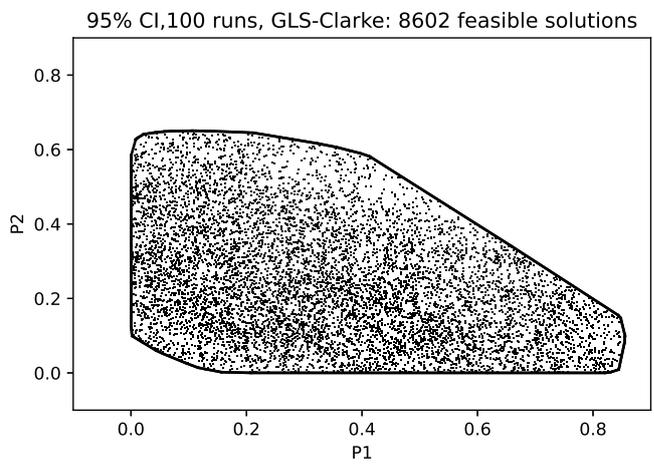


**Figure 6** Relationship between coefficients of variation ( $CV$ ) and the number of suspended sediment samples of the 95% confidence regions for channel banks ( $P_1$ ) and unpaved roads ( $P_2$ ) when sample sizes are reduced by sequence  $nY = \{24, 20, 16, 12, 8, 4, 2\}$  using OLS\_Clarke and GLS\_Clarke models over 100 repetitions.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



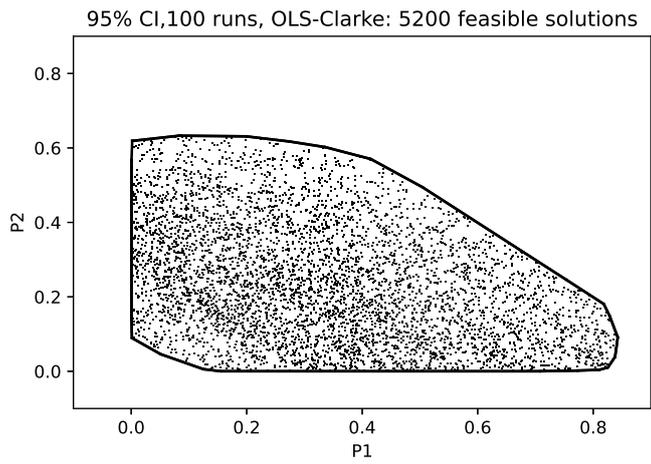
(a) OLS\_Clarke



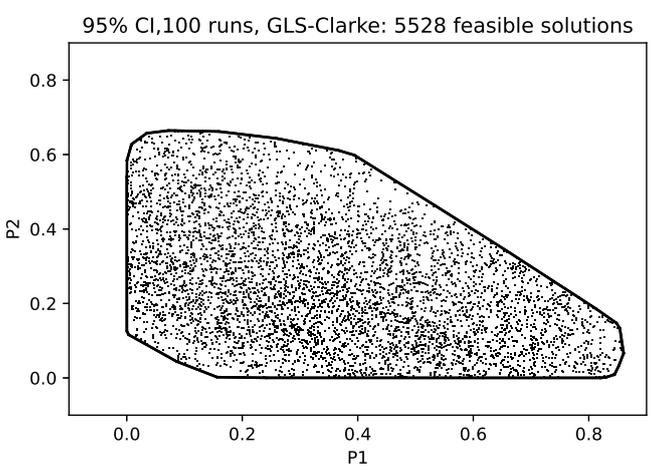
(b) GLS\_Clarke

**Figure 7** The feasible solutions in the 95% confidence region. The mean of proportions  $(\bar{P}_1, \bar{P}_2) = (0.341, 0.238)$  and  $(\bar{P}_1, \bar{P}_2) = (0.341, 0.245)$  obtained by a) OLS\_Clarke model and b) GLS\_Clarke model, respectively. Considering  $nL = 20$  over 100 repetitions.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



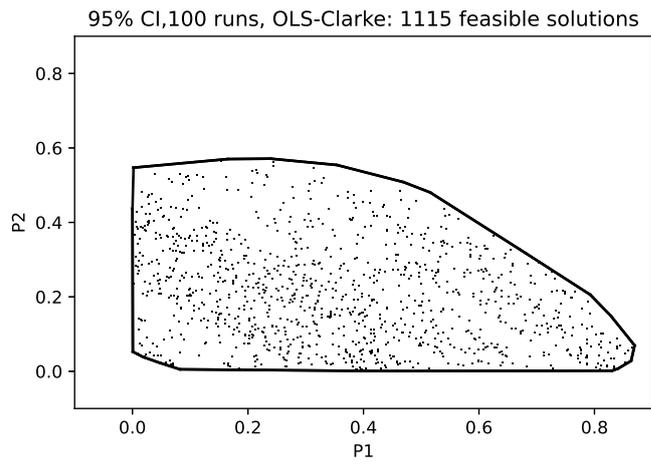
(a) OLS\_Clarke



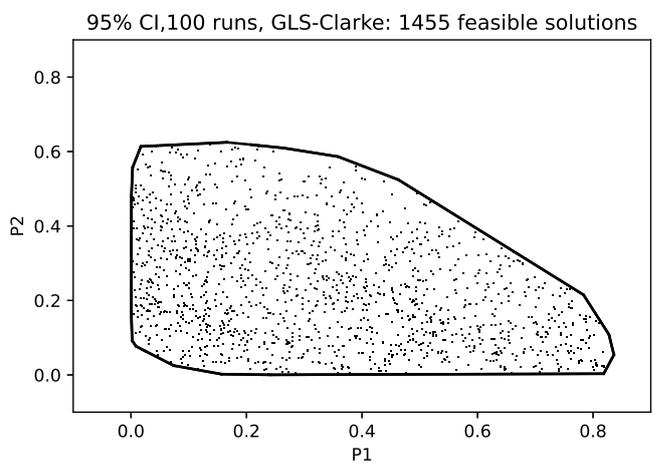
(b) GLS\_Clarke

**Figure 8** The feasible solutions in the 95% confidence region. The mean of proportions  $(\bar{P}_1, \bar{P}_2) = (0.338, 0.239)$  and  $(\bar{P}_1, \bar{P}_2) = (0.348, 0.252)$  obtained by a) OLS\_Clarke model and b) GLS\_Clarke model, respectively. Considering  $nL = 12$  over 100 repetitions.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



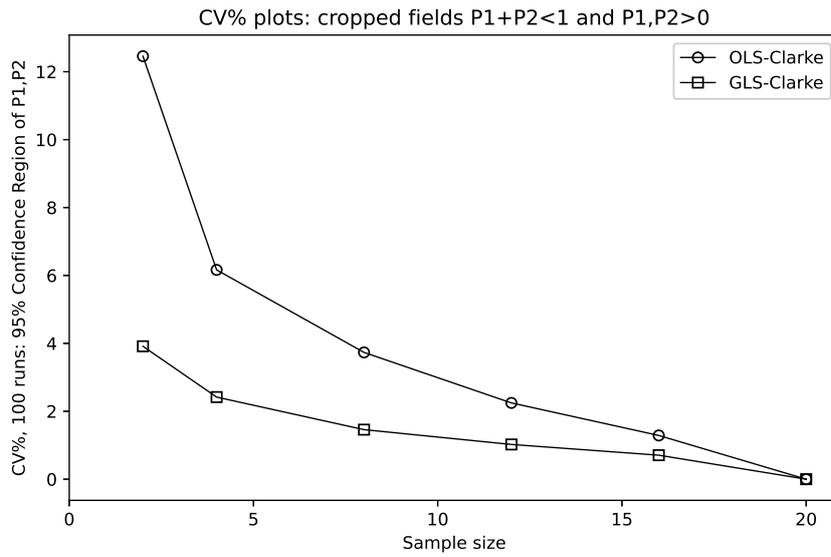
(a) OLS\_Clarke



(b) GLS\_Clarke

**Figure 9** The feasible solutions in the 95% confidence region. The mean of proportions  $(\bar{P}_1, \bar{P}_2) = (0.358, 0.208)$  and  $(\bar{P}_1, \bar{P}_2) = (0.330, 0.234)$  obtained by a) OLS\_Clarke model and b) GLS\_Clarke model, respectively. Considering  $nL = 4$  over 100 repetitions.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



**Figure 10** Relationship between coefficients of variation ( $CV$ ) and the number of suspended sediment samples of the 95% confidence regions for channel banks ( $P_1$ ) and unpaved roads ( $P_2$ ) when samples sizes are reduced by sequence  $nCF = \{20, 16, 12, 8, 4, 2\}$  using OLS\_Clarke and GLS\_Clarke models over 100 repetitions.

**Table 1** Simulation 1 - Considering the reduction of suspended sediment samples ( $nY$ ). a) OLS\_Clarke model and b) GLS\_Clarke model.

a) OLS_Clarke model							
nY	Mean Area	$\sigma$	$C_n$	$n^*$	$\overline{P_1}$	$\overline{P_2}$	$\overline{P_3}$
24	0.403	0.000	38880	7722	0.341	0.238	0.420
20	0.401	0.005	32400	6592	0.345	0.244	0.412
16	0.402	0.007	25920	5163	0.316	0.243	0.440
12	0.399	0.011	19440	4258	0.319	0.248	0.432
8	0.392	0.015	12960	2298	0.392	0.219	0.389
4	0.380	0.028	6480	1022	0.333	0.219	0.448
2	0.366	0.036	3240	779	0.325	0.318	0.356
b) GLS_Clarke model							
nY	Mean Area	$\sigma$	$C_n$	$n^*$	$\overline{P_1}$	$\overline{P_2}$	$\overline{P_3}$
24	0.418	0.000	38880	5550	0.259	0.276	0.466
20	0.418	0.004	32400	4447	0.263	0.286	0.451
16	0.418	0.005	25920	3435	0.263	0.281	0.457
12	0.415	0.007	19440	2850	0.267	0.278	0.456
8	0.413	0.010	12960	2240	0.256	0.270	0.474
4	0.401	0.016	6480	761	0.266	0.276	0.458
2	0.384	0.019	3240	409	0.282	0.318	0.400

Number of suspended sediment samples ( $nY$ ); Mean area of the 95% confidence regions; Standard deviations ( $\sigma$ ); Number of possible solutions of the overdetermined systems ( $C_n$ ); Feasible solutions ( $n^*$ ) where  $0 < P_1, P_2, P_3 < 1$  and  $P_1 + P_2 + P_3 = 1$ ; Means of proportions contributed by Channel Banks ( $\overline{P_1}$ ), Unpaved Roads ( $\overline{P_2}$ ) and Crop Fields ( $\overline{P_3}$ ) of the 95% confidence regions and over 100 repetitions.

**Table 2** Simulation 2 - Considering the reduction of crop fields samples ( $nL$ ). a) OLS\_Clarke model and b) GLS\_Clarke model.

a) OLS_Clarke model							
nL	Mean Area	$\sigma$	$C_n$	$n^*$	$\overline{P_1}$	$\overline{P_2}$	$\overline{P_3}$
20	0.403	0.000	38880	7722	0.341	0.238	0.420
16	0.401	0.005	31104	5934	0.352	0.241	0.407
12	0.400	0.009	23328	5200	0.338	0.239	0.423
8	0.394	0.015	15552	3154	0.350	0.232	0.418
4	0.386	0.024	7776	1115	0.358	0.208	0.434
2	0.351	0.044	3888	547	0.347	0.214	0.439
b) GLS_Clarke model							
nL	Mean Area	$\sigma$	$C_n$	$n^*$	$\overline{P_1}$	$\overline{P_2}$	$\overline{P_3}$
20	0.414	0.000	38880	8602	0.341	0.245	0.414
16	0.414	0.003	31104	6862	0.344	0.247	0.408
12	0.413	0.004	23328	5528	0.348	0.252	0.400
8	0.410	0.006	15552	3242	0.345	0.244	0.411
4	0.405	0.010	7776	1455	0.330	0.234	0.437
2	0.392	0.015	3888	684	0.319	0.232	0.449

Number of crop fields samples ( $nCF$ ); Mean area of the 95% confidence regions; Standard deviations ( $\sigma$ ); Number of possible solutions of the overdetermined systems ( $C_n$ ); Feasible solutions ( $n^*$ ) where  $0 < P_1, P_2, P_3 < 1$  and  $P_1 + P_2 + P_3 = 1$ ; Means of proportions contributed by Channel Banks ( $\overline{P_1}$ ), Unpaved Roads ( $\overline{P_2}$ ) and Crop Fields ( $\overline{P_3}$ ) of the 95% confidence regions and over 100 repetitions.

## References

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 641  
642 [1] Owens, P.: Soil erosion and sediment dynamics in the anthropocene: a re-  
643 view of human impacts during a period of rapid global environmental change.  
644 Journal of Soils and Sediments **20**, 1–29 (2020) <https://doi.org/10.1007/s11368-020-02815-9>  
645
- 646 [2] Collins, A.L., Blackwell, M., Boeckx, P., Chivers, C.-A., Emelko, M., Evrard, O.,  
647 Foster, I., Gellis, A., Gholami, H., Granger, S., *et al.*: Sediment source fingerprint-  
648 ing: benchmarking recent outputs, remaining challenges and emerging themes.  
649 Journal of Soils and Sediments **20**(12), 4160–4193 (2020)
- 650 [3] Minella, J.P.G., Walling, D.E., Merten, G.H.: Establishing a sediment budget for  
651 a small agricultural catchment in southern brazil, to support the development of  
652 effective sediment management strategies. Journal of Hydrology **519**, 2189–2201  
653 (2014) <https://doi.org/10.1016/j.jhydrol.2014.10.013>
- 654 [4] Walling, D.E.: The evolution of sediment source fingerprinting investigations in  
655 fluvial systems. Journal of Soils and Sediments **13**(10), 1658–1675 (2013)
- 656 [5] Laceyby, J.P., Olley, J.: An examination of geochemical modelling approaches to  
657 tracing sediment sources incorporating distribution mixing and elemental corre-  
658 lations. Hydrological Processes **29**(6), 1669–1685 (2015) <https://doi.org/10.1002/hyp.10287> <https://onlinelibrary.wiley.com/doi/pdf/10.1002/hyp.10287>
- 660 [6] Franks, S., Rowan, J.: Multi-parameter fingerprinting of sediment sources: Uncer-  
661 tainty estimation and tracer selection. Computational methods in water resources  
662 - Volume 2 - Computational methods,surface water systems and hydrology,  
663 1067–1074 (2000)
- 664 [7] Collins, A., Walling, D.: Selecting fingerprint properties for discriminating po-  
665 tential suspended sediment sources in river basins. Journal of Hydrology **261**,  
666 218–244 (2002)
- 667 [8] Cooper, R.J., Krueger, T., Hiscock, K.M., Rawlins, B.G.: Sensi-  
668 tivity of fluvial sediment source apportionment to mixing model  
669 assumptions: A bayesian model comparison. Water Resources Re-  
670 search **50**(11), 9031–9047 (2014) <https://doi.org/10.1002/2014WR016194>  
671 <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1002/2014WR016194>
- 672 [9] Pulley, S., Collins, A.L.: Tracing catchment fine sediment sources using the new  
673 sift (sediment fingerprinting tool) open source software. Science of The Total  
674 Environment **635**, 838–858 (2018) <https://doi.org/10.1016/j.scitotenv.2018.04.126>  
675
- 676 [10] Hughes, A.O., Olley, J.M., Croke, J.C., McKergow, L.A.: Sediment source

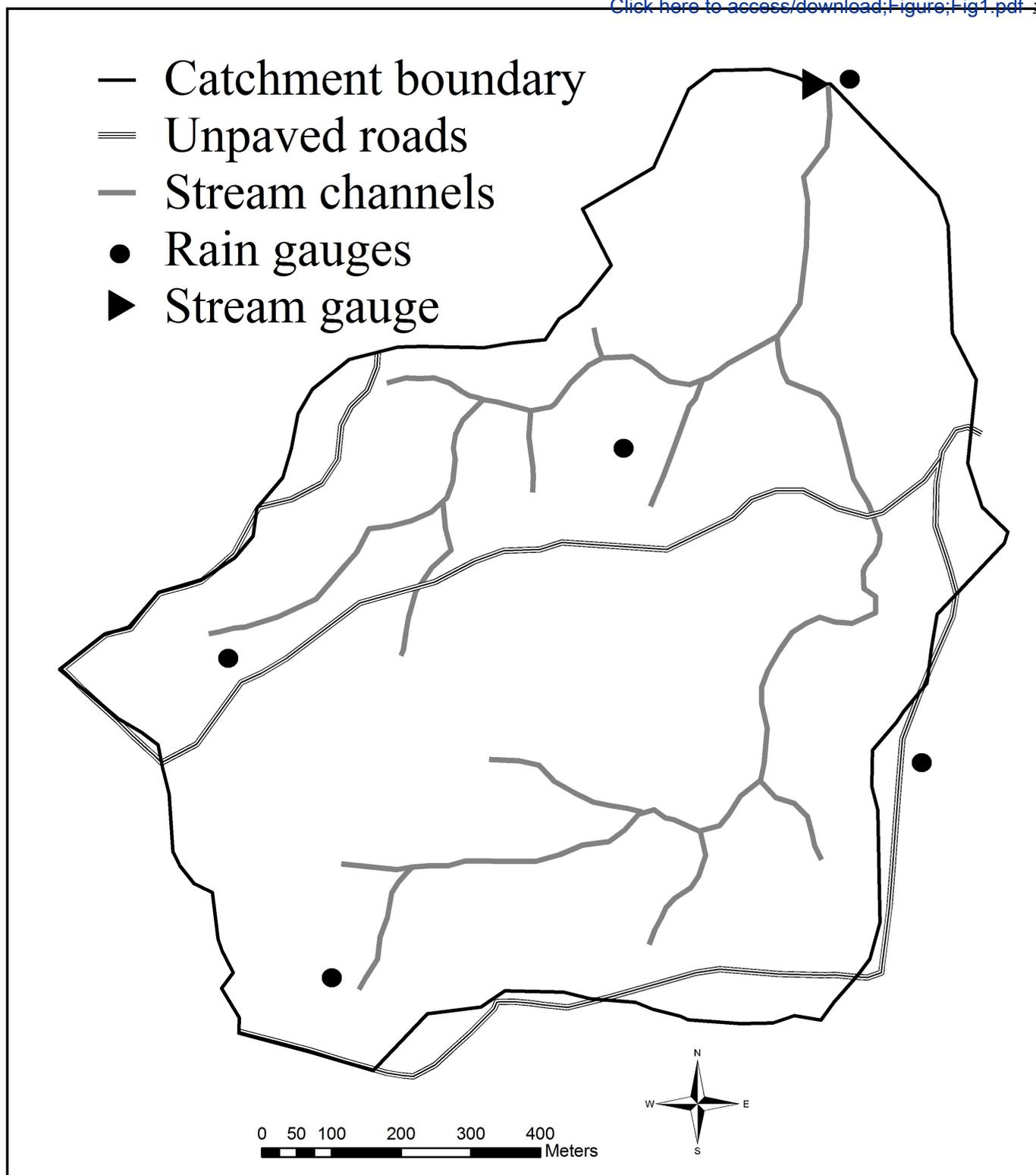
- 677 changes over the last 250 years in a dry-tropical catchment, central queens-  
678 land, australia. *Geomorphology* **104**(3), 262–275 (2009) [https://doi.org/10.1016/](https://doi.org/10.1016/j.geomorph.2008.09.003)  
679 [j.geomorph.2008.09.003](https://doi.org/10.1016/j.geomorph.2008.09.003)
- 680 [11] Uber, M., Legout, C., Nord, G., Crouzet, C., Demory, F., Poulenard, J.: Compar-  
681 ing alternative tracing measurements and mixing models to fingerprint suspended  
682 sediment sources in a mesoscale mediterranean catchment. *Journal of Soils and*  
683 *Sediments* **19**, 3255–3273 (2019)
- 684 [12] Lizaga, I., Latorre, B., Gaspar, L., Navas, A.: Consensus ranking as a method to  
685 identify non-conservative and dissenting tracers in fingerprinting studies. *Science*  
686 *of the Total Environment* **720**, 137537 (2020)
- 687 [13] Latorre, B., Lizaga, I., Gaspar, L., Navas, A.: A novel method for analysing  
688 consistency and unravelling multiple solutions in sediment fingerprinting. *Science*  
689 *of The Total Environment* **789**, 147804 (2021)
- 690 [14] Batista, P., Lacey, J., Evrard, O.: How to evaluate sediment fingerprinting source  
691 apportionments. *Journal of Soils and Sediments* **22**, 1–14 (2022) [https://doi.org/](https://doi.org/10.1007/s11368-022-03157-4)  
692 [10.1007/s11368-022-03157-4](https://doi.org/10.1007/s11368-022-03157-4)
- 693 [15] Collins, A., Pulley, S., Foster, I.D., Gellis, A., Porto, P., Horowitz, A.: Sediment  
694 source fingerprinting as an aid to catchment management: a review of the current  
695 state of knowledge and a methodological decision-tree for end-users. *Journal of*  
696 *Environmental Management* **194**, 86–108 (2017)
- 697 [16] Lacey, J., Gellis, A., Koiter, A., Blake, W., Evrard, O.: Preface—evaluating  
698 the response of critical zone processes to human impacts with sediment source  
699 fingerprinting. *Journal of Soils and Sediments* **19** (2019) [https://doi.org/10.1007/](https://doi.org/10.1007/s11368-019-02409-0)  
700 [s11368-019-02409-0](https://doi.org/10.1007/s11368-019-02409-0)
- 701 [17] Evrard, O., Batista, P.V., Company, J., Dabrin, A., Foucher, A., Frankl, A.,  
702 García-Comendador, J., Hugué, A., Lake, N., Lizaga, I., *et al.*: Improving the  
703 design and implementation of sediment fingerprinting studies: Summary and out-  
704 comes of the tracing 2021 scientific school. *Journal of Soils and Sediments* **22**(6),  
705 1648–1661 (2022)
- 706 [18] Walling, D., Woodward, J.: Tracing sources of suspended sediment in river basins:  
707 A case study of the river culm, devon, uk. *Marine and Freshwater Research* **46**(1),  
708 327–336 (1995) <https://doi.org/10.1071/MF9950327>
- 709 [19] Chatterjee, S., Hadi, A.S.: *Regression Analysis by Example*. John Wiley & Sons,  
710 Chichester (2013)
- 711 [20] Montgomery, D.C., Peck, E.A., Vining, G.G.: *Introduction to Linear Regression*  
712 *Analysis*. John Wiley & Sons, ??? (2021)

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 713 [21] Kariya, T., Kurata, H.: Generalized Least Squares. Wiley Series in Probability  
714 and Statistics, p. 312. John Wiley & Sons, Chichester (2004)
- 715 [22] Baksalary, J., Puntanen, S.: Weighted-least-squares estimation in the general  
716 gauss-markov model. In: Dodge, Y. (ed.) Statistical Data Analysis and Inference,  
717 pp. 355–368. Elsevier Science Publishers B.V., ??? (1989)
- 718 [23] Ciarlet, P.G., Lions, J.L.: Handbook of Numerical Analysis. Vol. 1: Finite Dif-  
719 ference Methods (Part 1) and Solution of Equations in  $\mathbb{R}^n$  (Part 1). Elsevier  
720 Science Publishers, North Holland (1990)
- 721 [24] Meyer, C.D.: Matrix Analysis and Applied Linear Algebra. Other Titles in  
722 Applied Mathematics, vol. 71, p. 730. Society for Industrial and Applied  
723 Mathematics (SIAM), Philadelphia (2000)
- 724 [25] Shores, T.S.: Applied Linear Algebra and Matrix Analysis, p. 479. Springer, Berlin  
725 (2007). <https://doi.org/10.1007/978-3-319-74748-4>
- 726 [26] Lizaga, I., Latorre, B., Gaspar, L., Navas, A.: FingerPro: an R Package for  
727 Tracking the Provenance of Sediment. Water Resources Management: An In-  
728 ternational Journal, Published for the European Water Resources Association  
729 (EWRA) **34**(12), 3879–3894 (2020) <https://doi.org/10.1007/s11269-020-02650>
- 730 [27] Golub, G.H., Van Loan, C.F.: Matrix Computations. Johns Hopkins University  
731 Press, Baltimore (2013)
- 732 [28] Gentle, J.E.: Matrix algebra. Springer texts in statistics, Springer, New York, NY  
733 **10**, 978 (2007)
- 734 [29] Lawson, C.L., Hanson, R.J.: Solving least squares problems. In: Classics in  
735 Applied Mathematics (1976)
- 736 [30] Stock, B.C., Jackson, A.L., Ward, E.J., Parnell, A.C., Phillips, D.L., Semmens,  
737 B.X.: Analyzing mixing systems using a new generation of bayesian tracer mixing  
738 models. PeerJ **6**, 5096 (2018)
- 739 [31] Pulley, S., Collins, A.: Tracing catchment fine sediment sources using the new  
740 sift (sediment fingerprinting tool) open source software. Science of The Total  
741 Environment **635** (2018) <https://doi.org/10.1016/j.scitotenv.2018.04.126>
- 742 [32] Haddadchi, A., Ryder, D.S., Evrard, O., Olley, J.: Sediment fingerprinting in flu-  
743 vial systems: review of tracers, sediment sources and mixing models. International  
744 Journal of Sediment Research **28**(4), 560–578 (2013)
- 745 [33] Haddadchi, A., Olley, J., Laceby, J.: Accuracy of mixing models in predicting  
746 sediment source contributions. The Science of the total environment **497-498C**,  
747 139–152 (2014) <https://doi.org/10.1016/j.scitotenv.2014.07.105>

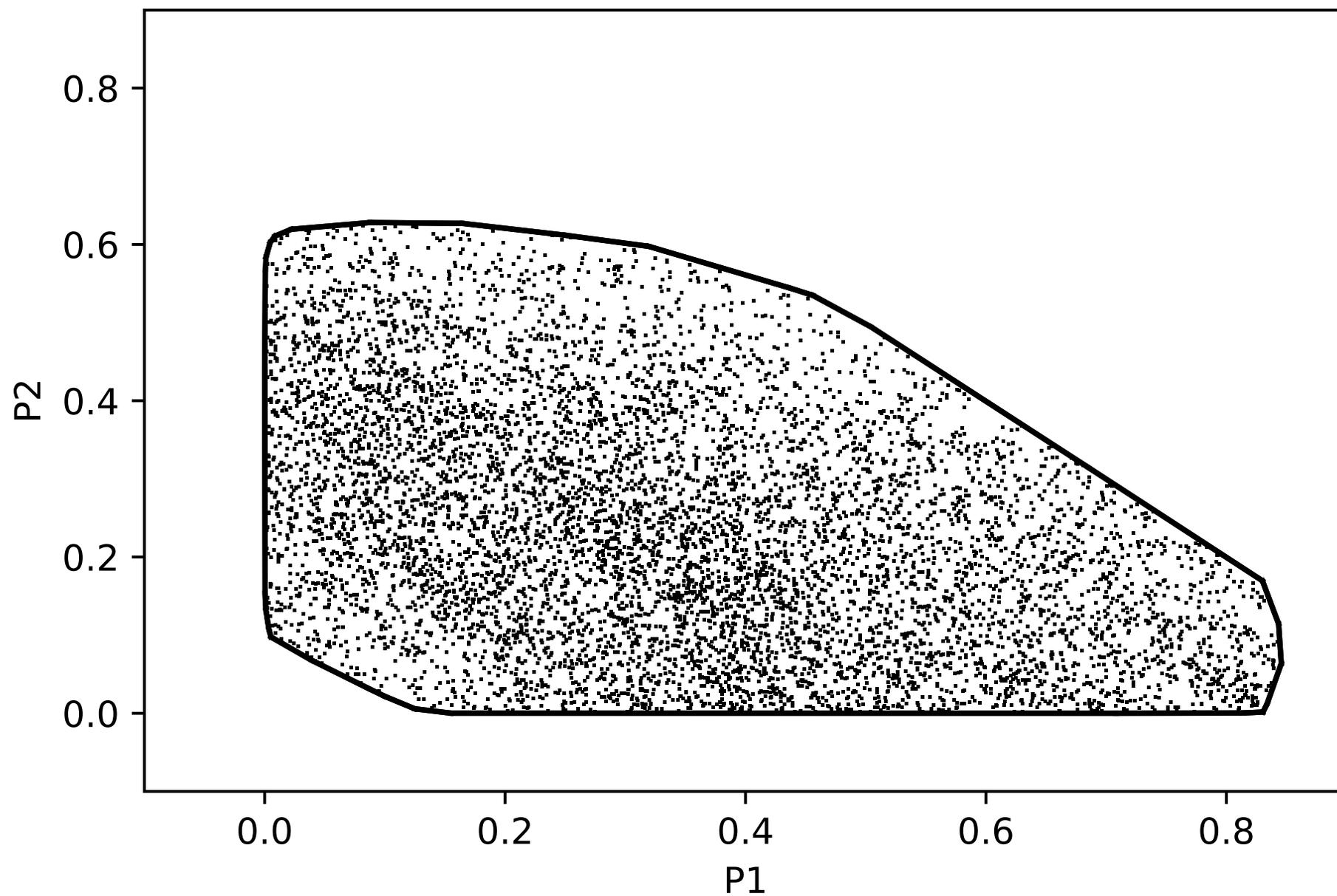
- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 748 [34] Rowan, J., Goodwill, P., Franks, S.: Uncertainty estimation in fingerprinting  
749 suspended sediment sources. In: Tracers in Geomorphology, pp. 279–290 (2000)
- 750 [35] Collins, A., Walling, D.: Sources of fine sediment recovered from the channel bed  
751 of lowland groundwater fed catchments in the uk. *Geomorphology* **88**, 120–138  
752 (2007) <https://doi.org/10.1016/j.geomorph.2006.10.018>
- 753 [36] Clarke, R.T.: A bootstrap calculation of confidence regions for proportions  
754 of sediment contributed by different source areas in a ‘fingerprinting’ model.  
755 *Hydrological Processes* **29**(12), 2694–2703 (2015) [https://doi.org/10.1002/hyp.  
756 10397](https://doi.org/10.1002/hyp.10397)
- 757 [37] Collins, A.L., Zhang, Y., Walling, D.E., Grenfell, S.E., Smith, P., Grischeff, J.,  
758 Locke, A., Sweetapple, A., Brogden, D.: Quantifying fine-grained sediment sources  
759 in the river axe catchment, southwest england: application of a monte carlo numer-  
760 ical modelling framework incorporating local and genetic algorithm optimisation.  
761 *Hydrological Processes* **26**(13), 1962–1983 (2012) [https://doi.org/10.1002/hyp.  
762 8283](https://doi.org/10.1002/hyp.8283) <https://onlinelibrary.wiley.com/doi/pdf/10.1002/hyp.8283>
- 763 [38] Clarke, R.T., Minella, J.P.G.: Evaluating sampling efficiency when estimating  
764 sediment source contributions to suspended sediment in rivers by fingerprinting.  
765 *Hydrological Processes* **30**(19), 3408–3419 (2016) [https://doi.org/10.1002/hyp.  
766 10866](https://doi.org/10.1002/hyp.10866)
- 767 [39] Batista, P.V., Laceby, J.P., Davies, J., Carvalho, T.S., Tassinari, D., Silva, M.L.,  
768 Curi, N., Quinton, J.N.: A framework for testing large-scale distributed soil ero-  
769 sion and sediment delivery models: Dealing with uncertainty in models and the  
770 observational data. *Environmental Modelling & Software* **137**, 104961 (2021)
- 771 [40] Yu, L., Oldfield, F.: A multivariate mixing model for identifying sediment source  
772 from magnetic measurements. *Quaternary Research* **32**(2), 168–181 (1989) [https://doi.org/10.1016/0033-5894\(89\)90073-2](https://doi.org/10.1016/0033-5894(89)90073-2)
- 774 [41] De Maesschalck, R., Jouan-Rimbaud, D., Massart, D.L.: The mahalanobis  
775 distance. *Chemometrics and intelligent laboratory systems* **50**(1), 1–18 (2000)
- 776 [42] Daszykowski, M., Kaczmarek, K., Heyden, Y.V., Walczak, B.: Robust statis-  
777 tics in data analysis — a review: Basic concepts. *Chemometrics and Intelligent  
778 Laboratory Systems* **85**, 203–219 (2007)
- 779 [43] Minella, J.P., Walling, D.E., Merten, G.H.: Combining sediment source tracing  
780 techniques with traditional monitoring to assess the impact of improved land  
781 management on catchment sediment yields. *Journal of Hydrology* **348**(3-4), 546–  
782 563 (2008)
- 783 [44] Minella, J.P., Merten, G.H., Schlesner, A., Bernardi, F., Barros, C.A., Tiecher,

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 784 T., Ramon, R., Evrard, O., Santos, D.R., Reichert, J.M., *et al.*: Combining sed-  
785 iment source tracing techniques with traditional monitoring: The “arvorezinha  
786 catchment” experience. *Hydrological Processes* **36**(9), 14665 (2022)
- 787 [45] Davis, J.C., Sampson, R.J.: *Statistics and Data Analysis in Geology*, p. 646. Wiley,  
788 New York (1986)
- 789 [46] Harris, C., Millman, K., Walt, S., Gommers, R., Virtanen, P., Cournapeau, D.,  
790 Wieser, E., Taylor, J., Berg, S.: Smith 474 nj. Kern R, Picus M, Hoyer S, van  
791 Kerkwijk MH, Brett M, Haldane A, del R'io JF, Wiebe M, Peterson P, G'erard-  
792 475 Marchant P, et al. Array programming with NumPy. *Nature* **585**(7825), 357–  
793 362 (2020) <https://doi.org/10.1038/s41586-020-2649-2>
- 794 [47] Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cour-  
795 napeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., *et al.*: Scipy  
796 1.0: fundamental algorithms for scientific computing in python. *Nature methods*  
797 **17**(3), 261–272 (2020) <https://doi.org/10.1038/s41592-019-0686-2>
- 798 [48] Hunter, J.D.: Matplotlib: A 2d graphics environment. *Computing in science &*  
799 *engineering* **9**(03), 90–95 (2007) <https://doi.org/10.1109/MCSE.2007.55>
- 800 [49] Silva, C., Minella, J., Schlesner, A., Merten, G., Barros, C., Tassi, R., Dambroz,  
801 A.: Unpaved road conservation planning at the catchment scale. *Environmental*  
802 *monitoring and assessment* **193**(9), 595 (2021)
- 803 [50] Johnson, R.A., Wichern, D.W., *et al.*: *Applied Multivariate Statistical Analysis*.  
804 Pearson Education Inc, New Jersey (2002)

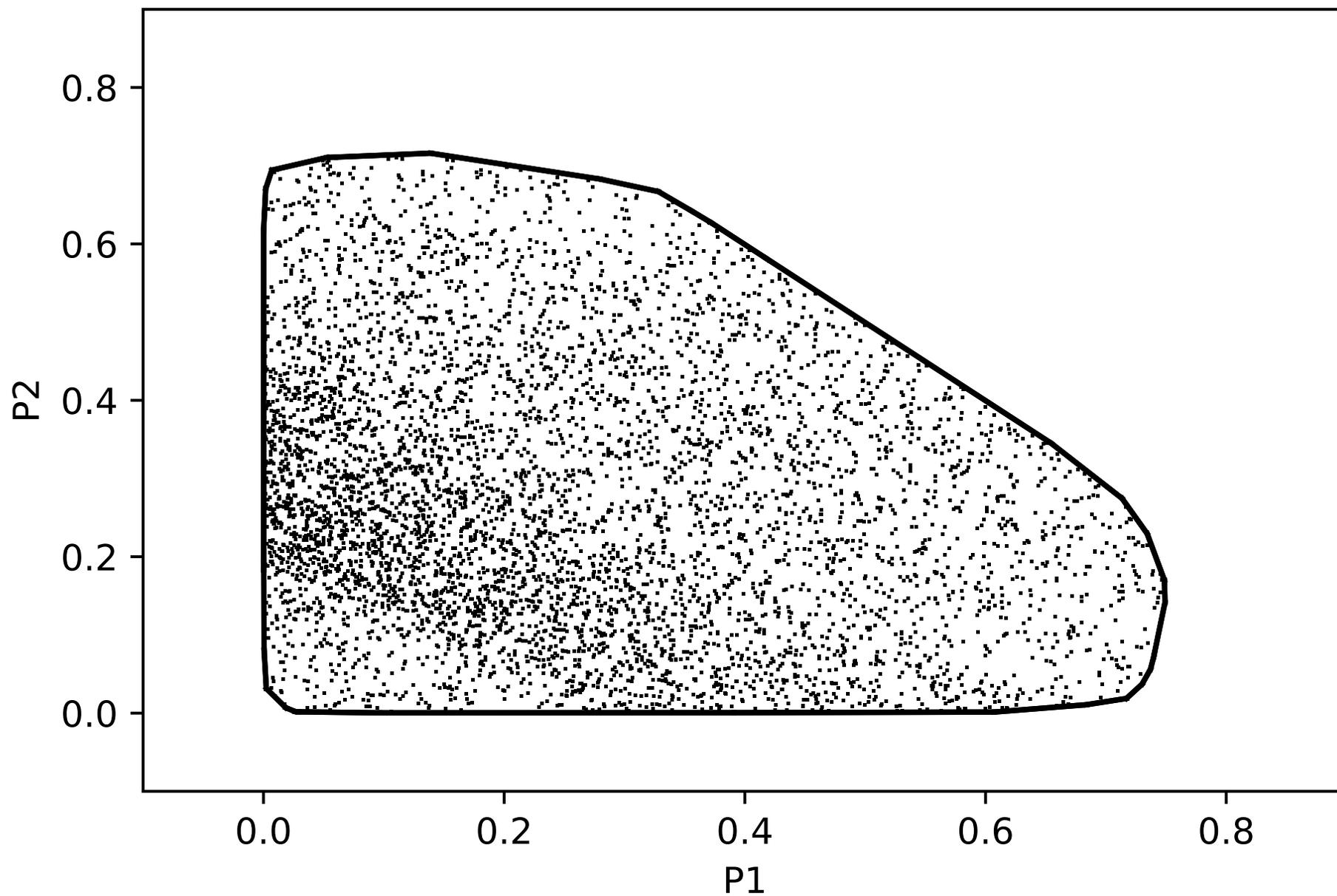
Figure 1



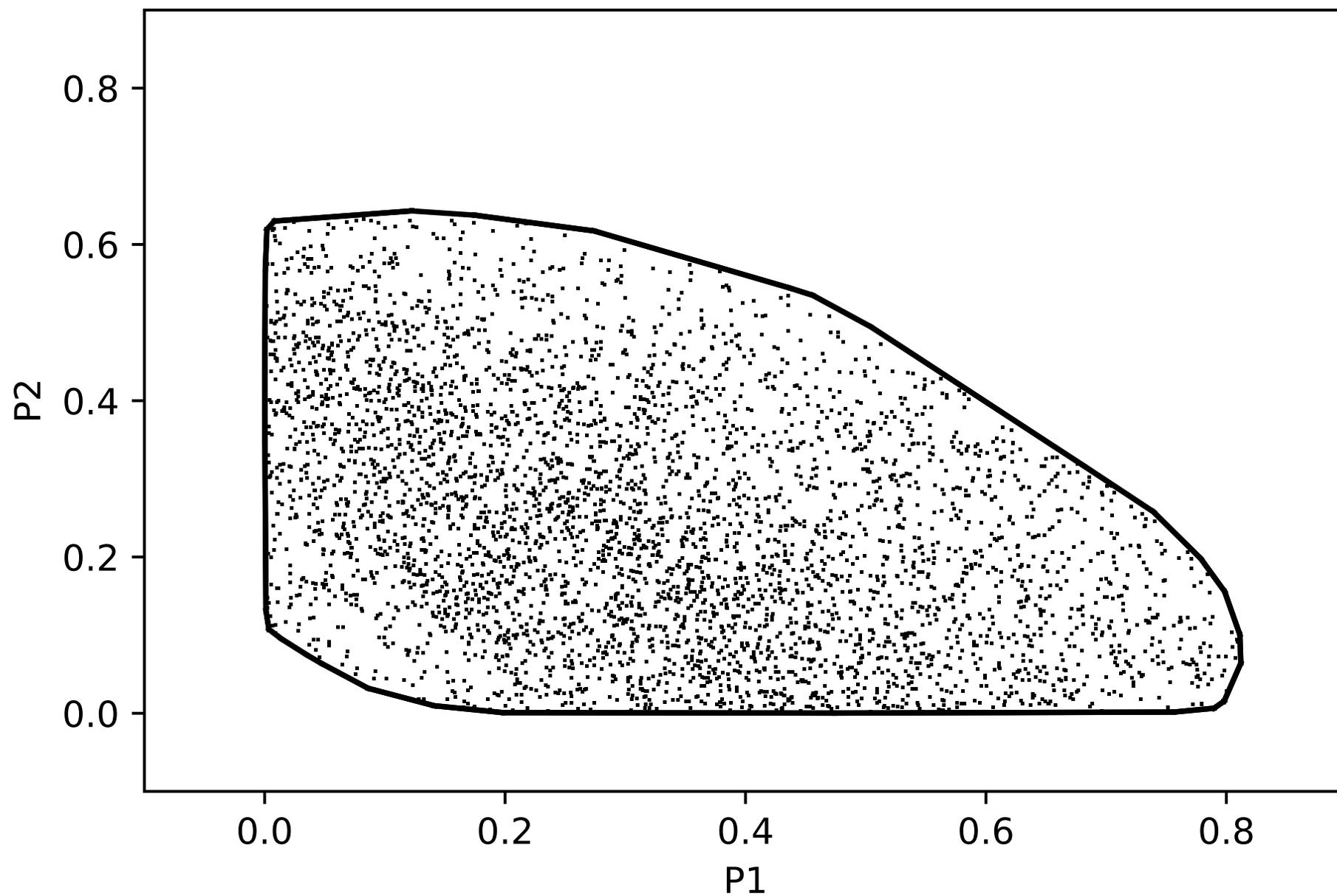
95% CI, 100 runs, OLS-Clarke: 7722 feasible solutions



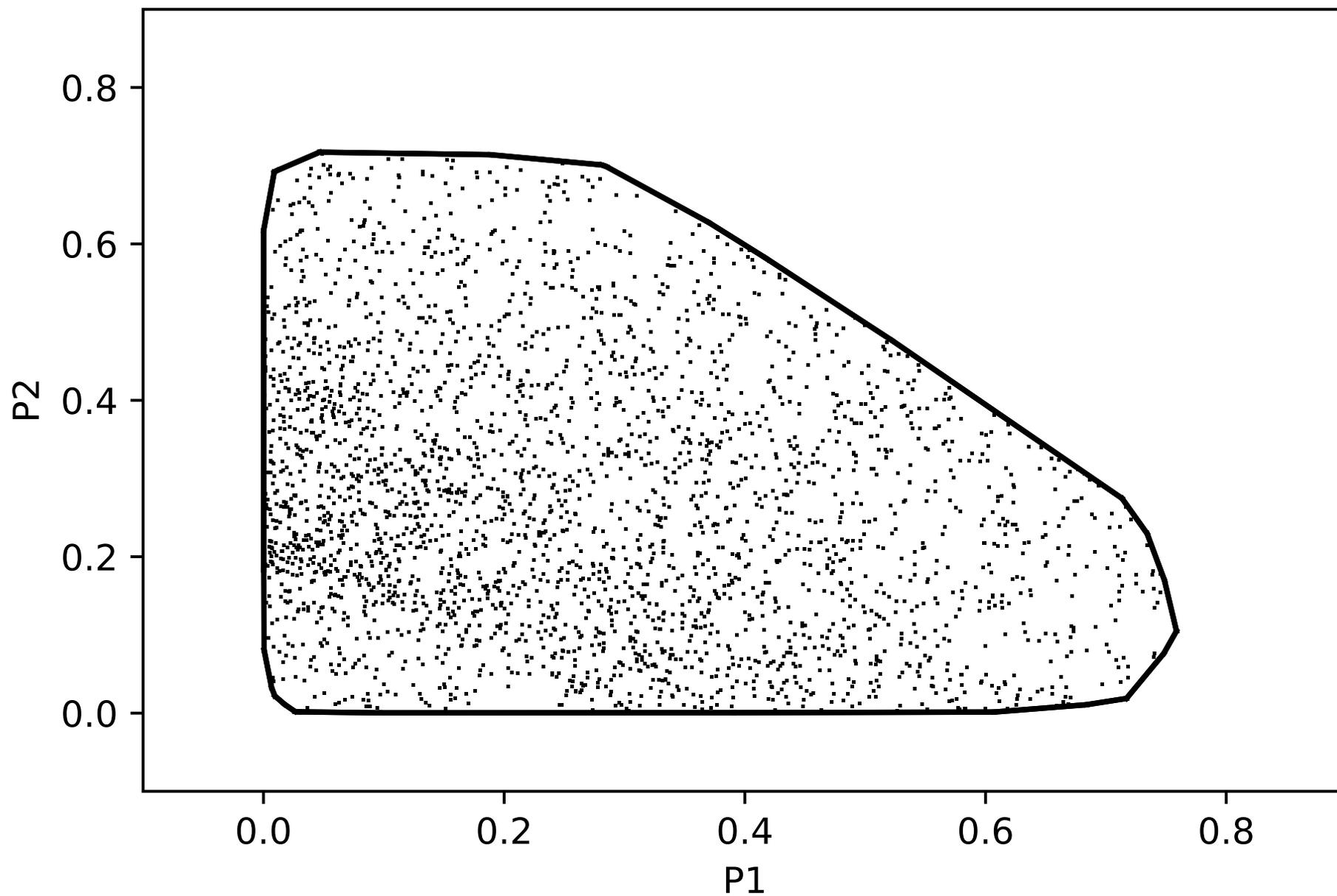
95% CI, 100 runs, GLS-Clarke: 5550 feasible solutions



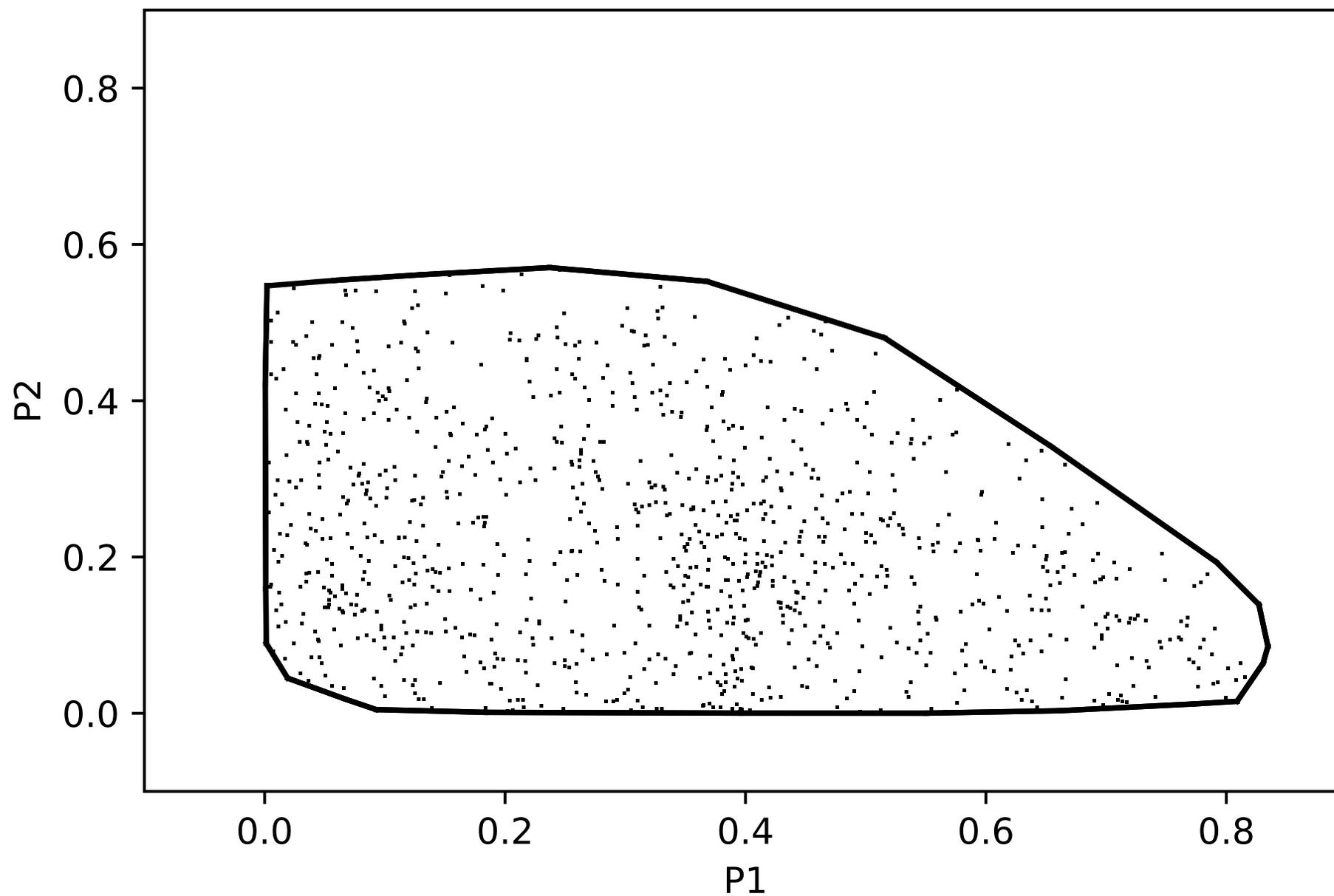
95% CI, 100 runs, OLS-Clarke: 4258 feasible solutions



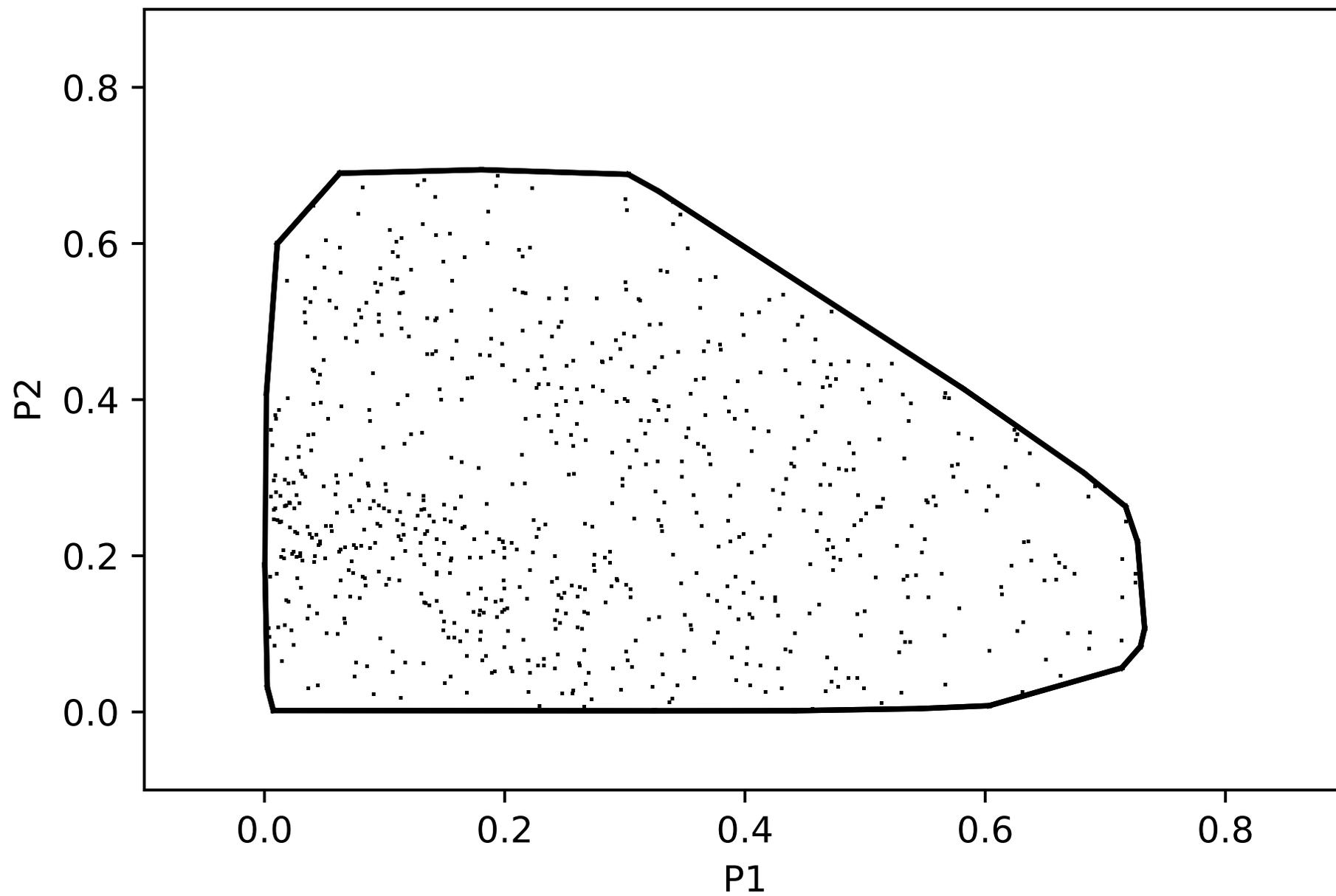
95% CI, 100 runs, GLS-Clarke: 2850 feasible solutions

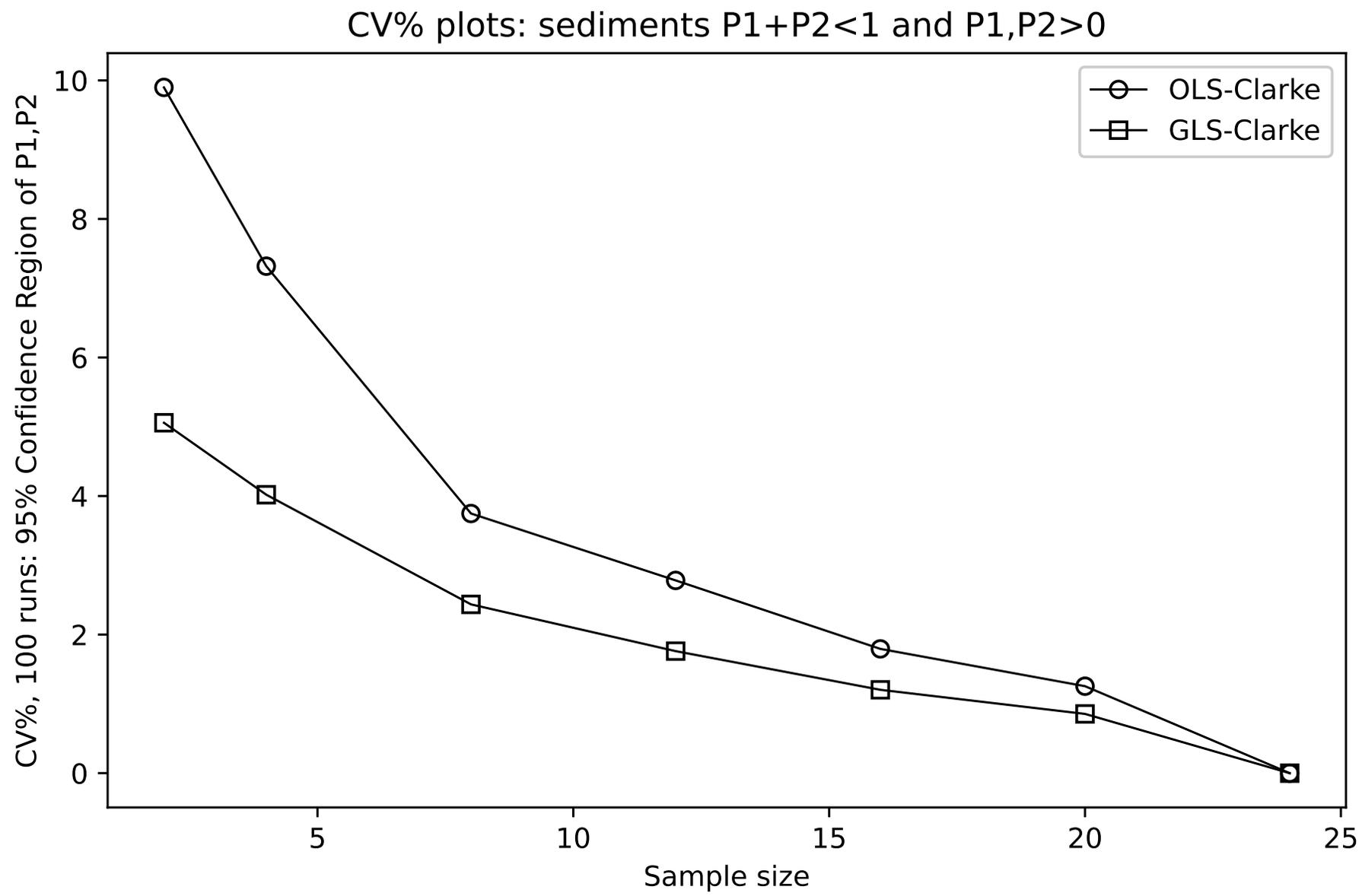


95% CI, 100 runs, OLS-Clarke: 1022 feasible solutions

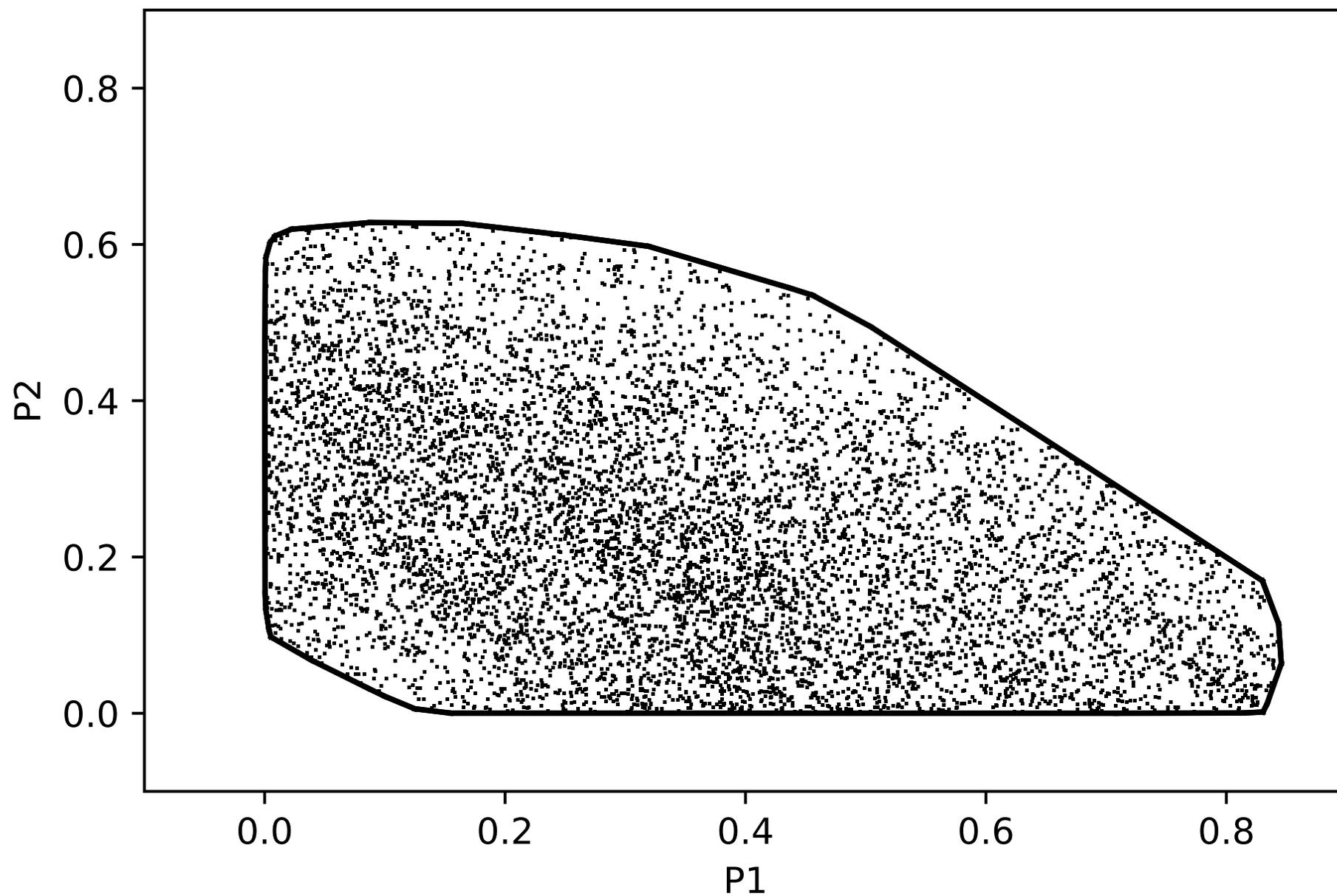


95% CI, 100 runs, GLS-Clarke: 761 feasible solutions

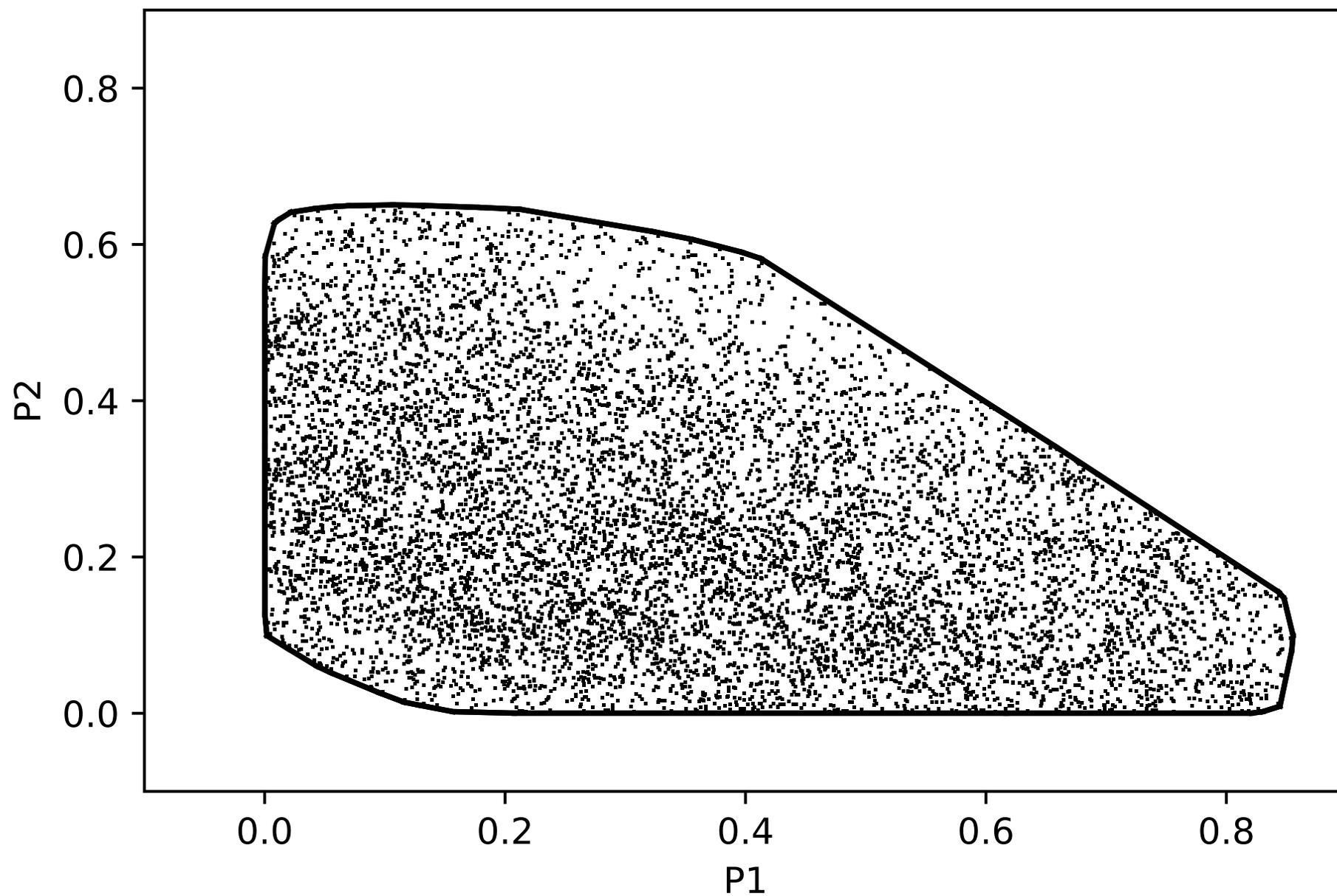




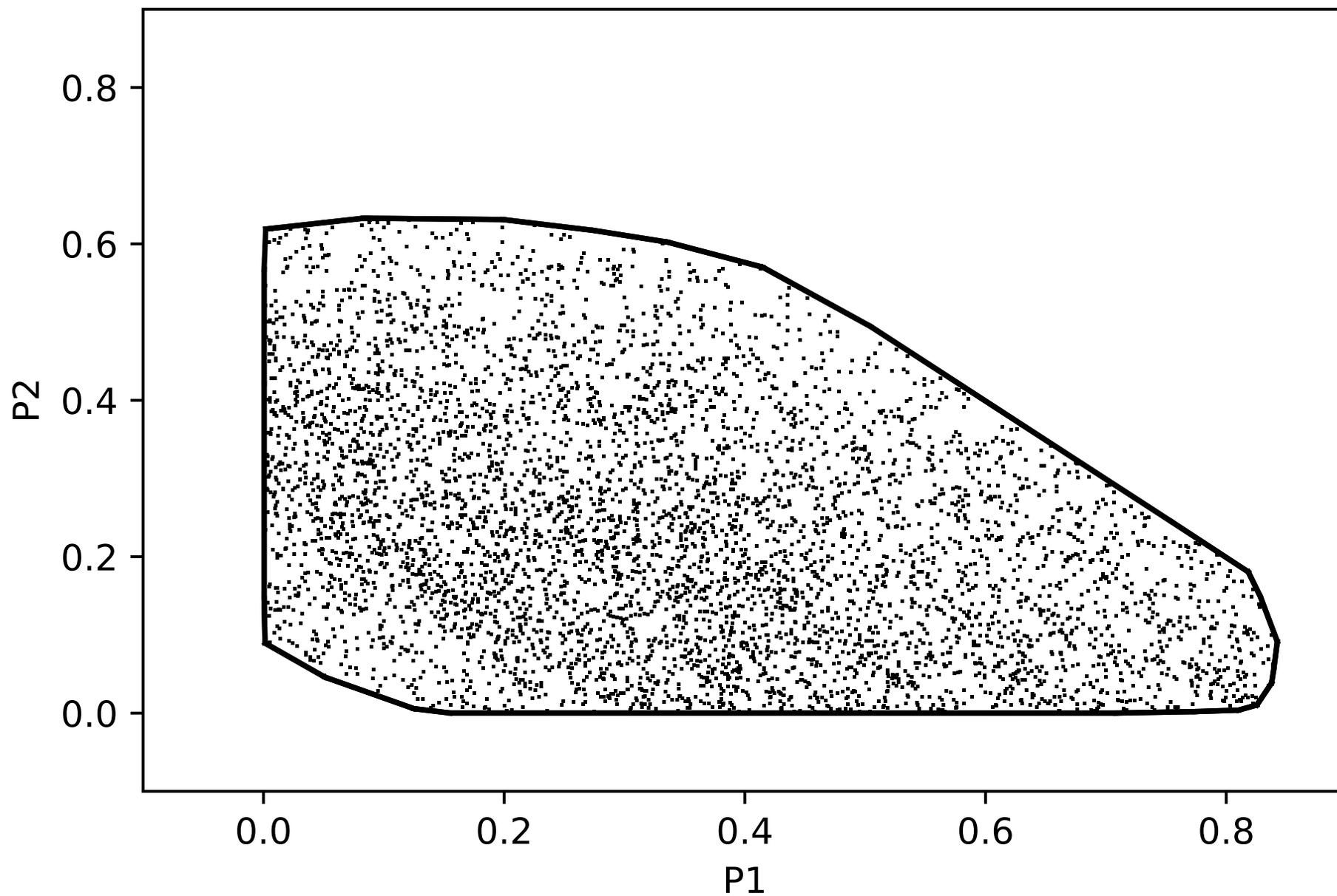
95% CI, 100 runs, OLS-Clarke: 7722 feasible solutions



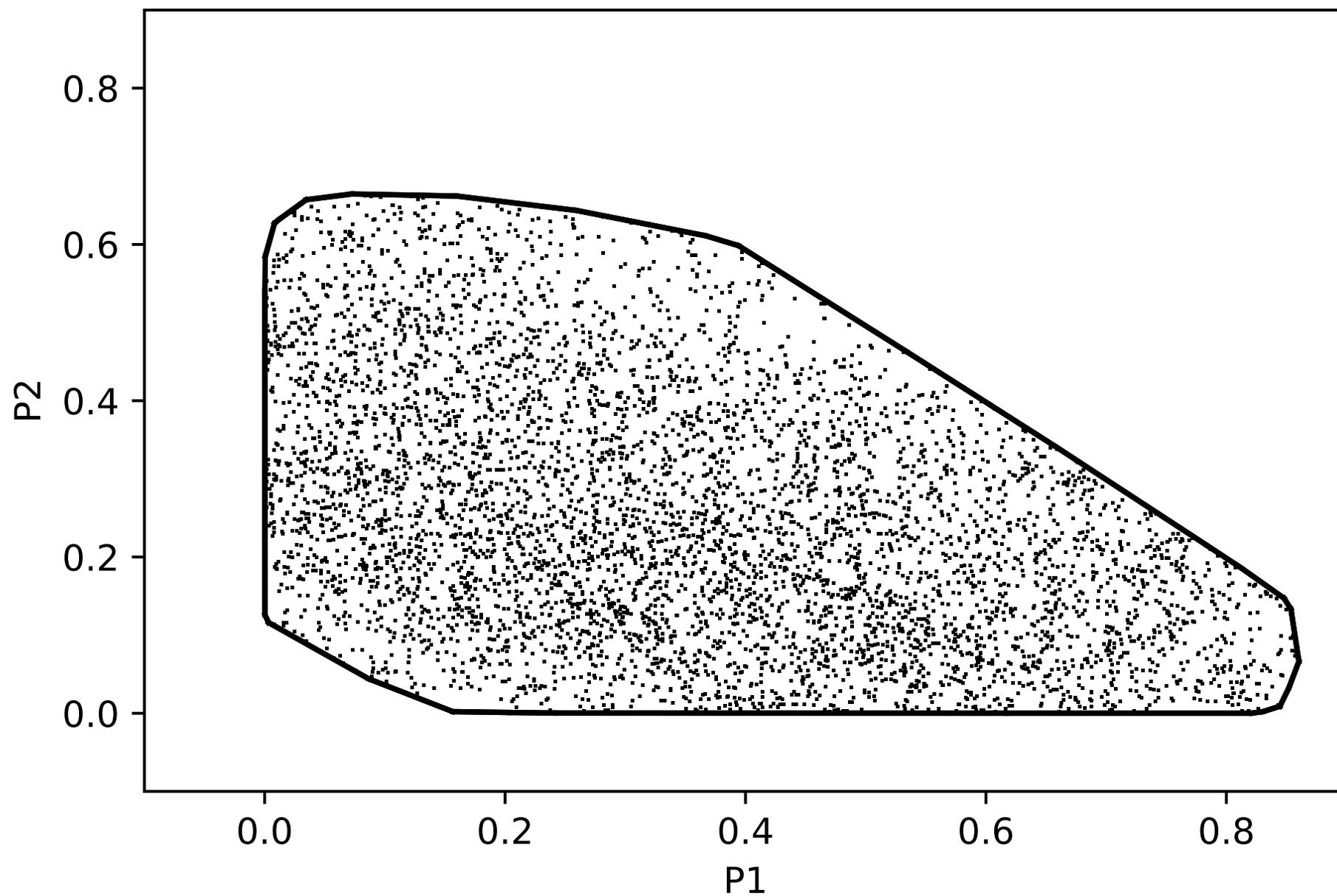
95% CI, 100 runs, GLS-Clarke: 8602 feasible solutions



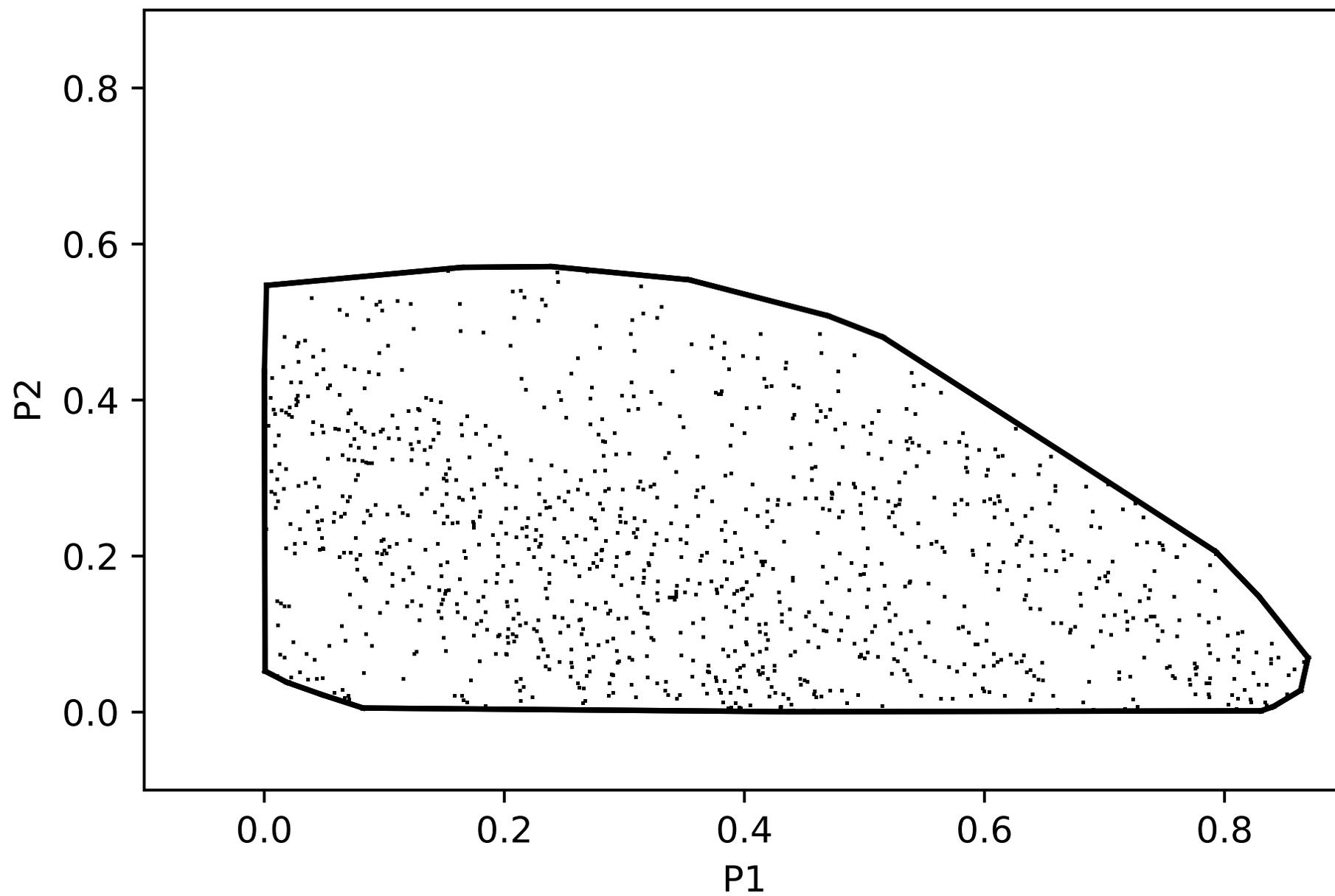
95% CI, 100 runs, OLS-Clarke: 5200 feasible solutions



95% CI, 100 runs, GLS-Clarke: 5528 feasible solutions



95% CI, 100 runs, OLS-Clarke: 1115 feasible solutions



95% CI, 100 runs, GLS-Clarke: 1455 feasible solutions

