1 **A machine learning-based approach for mapping leachate contamination**

2 **using geoelectrical methods**

3 Ester Piegari[1*], Giorgio De Donno[2], Davide Melegari[2], Valeria Paoletti[1]

4 [1] Dipartimento di Scienze della Terra, dell'Ambiente e delle Risorse, Università degli Studi di Napoli

5 Federico II, Naples, Italy

6 [2] Dipartimento di Ingegneria Civile Edile e Ambientale, "Sapienza" Università di Roma, Rome, Italy

7 * corresponding author: ester.piegari@unina.it

8

## 9 Abstract

10 The growth of urbanized areas combined with the overall increase of world's population is

11 leading to an increase of waste disposal sites that put a serious threat to the environment and

12 human health. Leachate is the main source of pollution in landfills and its negative impacts

13 continue for several years even after landfill closure. In recent years, geophysical methods are

14 recognized as effective tools for providing an imaging of the leachate plume. However, they

15 produce subsurface cross-sections in terms of individual physical quantities, leaving room for

16 ambiguities on interpretation of geophysical models and uncertainties in the definition of

17 contaminated zones. In this work, we propose a machine learning based approach for mapping

18 leachate contamination through an effective integration of geoelectrical tomographic data. We

19 apply the proposed approach for the characterization of two urban landfills. For both cases, we

20 use machine learning techniques to perform a multivariate analysis on datasets consisting of

21 electrical resistivity, chargeability and normalized chargeability (chargeability-to-resistivity

22 ratio) data extracted from previously inverted model sections. By executing a K-means

23 clustering analysis, we find the best partitioning of the datasets into different classes and get

24 updated cross-sections that provide a quantitative integration of the tomographic data, allowing

25 an objective identification of the most polluted zones. Our findings, also supported by borehole

26 data for one of the investigation sites, show that the combined use of geophysical imaging and

27 unsupervised machine learning is promising in environmental applications and can yield new

28 perspectives for the characterization of leachate distribution in landfills.

29

32

## 1. Introduction

34 Although there is an increased awareness on the importance of environment protection, urban

35 waste management is one the most important environmental issues. In many countries there is

36 still a little use of recycling or reuse actions, and the majority of municipal solid waste is

37 destined for landfills (WHO, 2015). Waste decomposition generates leachate that is a highly

38 contaminated liquid consisting of a mixture from organic degradation products, liquid waste

39 and rainwater. Leachate infiltration causes serious environmental issues to groundwater and

40 soils and, therefore, identifying and monitoring its flow pathways has major implications on

41 designing a risk mitigation strategy (Mukherjee et al., 2015; Lavagnolo et al., 2019; Vaccari et

42 al., 2019; Morita et al., 2021; Ergene et al., 2022). To this aim, geophysical methods often

43 represent the only cost-effective, rapid and non-invasive choice for mapping large areas, such

44 as those encountered in urban landfills, down to tens of meters (e.g., Di Maio et al., 2018).

45 In last decades, many studies have demonstrated that geoelectrical methods can be effective in

46 identifying landfill leachate (e.g., Soupios et al., 2007; De Donno and Cardarelli, 2017; Raji

47 and Adeoye, 2017; Soupios et al., 2017; Power et al., 2018, Zaini et al., 2022). As leachate is

48 a fluid with high concentrations of ions, it can be successfully imaged by low values of

49 electrical resistivity $\rho$ and high values of chargeability $M$, sensed respectively by electrical

50 resistivity tomography (ERT) and induced polarization (IP) surveys (Everett, 2013). The

51 contrast between the electrical properties of leachate and those of the surrounding media

52 generally facilitates its identification. However, geophysical inverted models leave ambiguities

53 in defining contaminated zones, particularly in presence of clayey soils, which are often placed

54 at the bottom of the landfill as a low-permeability barrier in combination with a synthetic liner

55 (geomembrane). Therefore, the zones characterized by the highest values of chargeability

56 frequently do not match with the most conductive ones, and the question of how to combine

57 information from different geophysical methods is still open. Many studies consider the so-

58 called normalized chargeability, $M_n$ (ratio of $M$ and $\rho$), as it is directly linked to surface

59 polarization (Slater and Lesmes, 2002). High values of $M_n$ are generally related to leachate

60 contamination, even though $M_n$ can be also largely and significantly affected by clay content

61 (Slater and Lesmes, 2002). Additionally, both resistivity and chargeability are related to the

62 saturation levels, and to pick up the threshold of $M_n$ for fully saturated zones (the most

63 dangerous ones) remains a subjective choice. In fact, focusing only on the highest values of $M_n$

64 may lead to false-positive results in the identification of leachate, as they may not arise from

65 the concomitance of low $\rho$ and high $M$ values. For instance, large values of normalized

66 chargeability may be caused by only large $M$ related to the presence of pockets of clays with

67 high-chargeability or by extremely low resistivity values of leachate-saturated zones. Thus, a

68 residual uncertainty persists in properly identifying leachate accumulation zones.

69 With this study, we propose the application of machine learning techniques to combine $\rho$, $M$

70 and $M_n$ data with the aim to get one comprehensive section integrating all information from the

71 electrical models.

72 Machine Learning (ML) is a branch of artificial intelligence based on the idea that systems can

73 learn from data, and this learning comes from analyses of data through statistical tools to make

74 predictions or finding patterns in data (Zhang et al., 2022). In recent years, applications of

75 clustering algorithms to geosciences are continuously growing due to their potential of

3

76  classifying large datasets into groups of data that share similar features (Lyra et al., 2014;

77  Lindsey et al., 2018; Bernardetti and Bruno, 2019; Karpatne et al., 2019; Straus, 2019; Abdideh

78  and Ameri, 2020, Kamer et al., 2020; Cesca. 2020; Piegari et al., 2022). There are many

79  different types of clustering algorithms, which can be broadly categorised into three types: i)

80  partitioning algorithms, which divide the dataset into a number of groups (clusters) and require

81  the number of clusters as an input data; ii) hierarchical algorithms, which provide a tree-based

82  representation of data points (dendrogram) showing the hierarchical relationship between

83  clusters; iii) density-based algorithms, which groups data points on the basis of their spatial

84  density. In this study, we use the K-means algorithm, which is one of the most common

85  partitioning algorithms. In addition to being a fast, robust and simple iterative algorithm, it has

86  the advantage of producing tighter (spherical) clusters than hierarchical and density-based

87  clustering (Shukla and Naganna, 2014; Zhang et al., 2022).

88  In the following sections, we describe the proposed methodology and show the results of

89  geophysical surveys performed into two urban waste landfills. Then, we report and discuss the

90  results of the K-means clustering analysis for the two case studies.

91

## 2. Methods

93  We illustrate our methodological approach in Fig. 1. Traditionally, to characterize leachate

94  contamination by means of geophysical methods, ERT and IP surveys are carried out with the

95  aim of retrieving two cross sections showing the distribution of $\rho$ and $M$, after data inversion.

96  In this work, we propose to go beyond traditional approaches by: (i) combining all inverted

97  values in a proper joint space and (ii) applying a cluster analysis to the unique large dataset.

98  The retrieved cluster indices are used to get integrated cross sections that allow direct

99  quantitative identification of leachate accumulation zones. In the following subsections, we

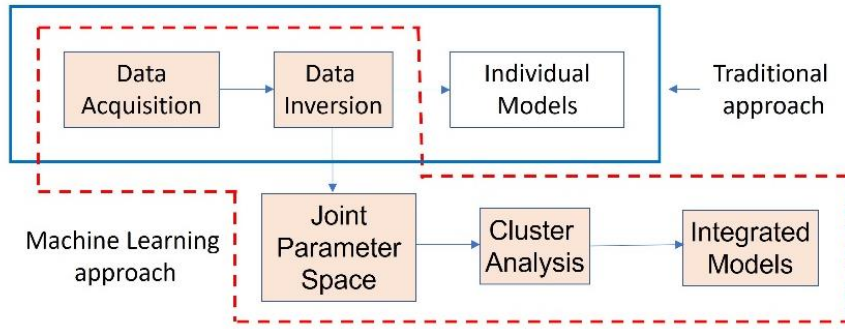100  give details about the inversion procedure and the clustering algorithms.

4

101

102   Figure 1. *Sketch of the traditional and ML based approaches. The red dotted line encloses the*

103   *proposed procedure.*

104

## 2.1   Geophysical Imaging

### 2.1.1 Forward modelling

107   The resistive response of a 2.5D subsoil (where the conductivity varies only in the *x-z* plane),

108   is described within a domain $D$ by the Fourier-transformed Poisson's equation under the

109   hypothesis of an external point source located at ($x_s$, $z_s$) (Dey and Morrison, 1979):

110   $$-\nabla \cdot [\sigma(x,z)\nabla\phi(x,z,\lambda)] + \lambda^2\sigma(x,z)\phi(x,z,\lambda) = I\delta(x_S)\delta(z_S) \quad \forall(x,z) \in D \,, \qquad (1)$$

111   where $\sigma$ is the conductivity ($\rho=1/\sigma$ is the resistivity), $\phi$ the transformed electric potential, $\lambda$

112   the transformed variable and $I$ the injected current.

113   Eq. (1), subjected to Dirichlet's and Neumann's boundary conditions on surface and lateral and

114   bottom boundaries respectively, is solved numerically. A widespread used technique is the

115   Galërkin formulation of the Finite Element Method (De Donno and Cardarelli, 2017), which is

116   employed in this work.

117   Once the solution of eq. (1) is achieved, potential is back-transformed and the apparent

118   resistivity can be predicted as:

119   $$\rho_a{}^{pre} = C\frac{\Delta V}{I} \,, \qquad (2)$$

5

where $C$ is the geometric factor and $\Delta V$ the potential difference, both obtained depending on the specific quadrupole sequence.

The capacitive response of a medium can be assessed in the time-domain through the chargeability $M$, which is proportional to the induced polarization (IP) phenomena occurring in subsoil due to the switch-on or switch-off of an external direct current (DC) source (Siegel, 1959).

In this case, the IP forward solution is given sequentially with the resistivity modelling by calculating the potential $V_M$ resulting from solution of eq. (1) but with the conductivity replaced by $\alpha = \sigma(1 - M)$.

The predicted apparent chargeability $M_a$ is then found as (Oldenbug and Li, 1994):

$$M_a{}^{pre} = \frac{V_M - V}{V_M}. \tag{3}$$

**2.1.2 Data inversion**

Time-domain measurements are performed employing a DC electrical source, measuring the apparent resistivity $\rho_a$ and apparent integral chargeability $M_a$, resulting from the potential decay after current switch-off (Binley and Slater, 2020):

$$\rho_a{}^{obs} = K \frac{V_p}{I}; \tag{4a}$$

$$M_a{}^{obs} = \frac{\int_{ti}^{tf} V_i dt}{V_p \Delta t}, \tag{4b}$$

where $V_p$ is the measured voltage during application of the DC current $I$ and $V_i$ the residual voltage after switch-off the electrical current, integrated over a time window $\Delta t$ defined between times $t_i$ and $t_f$. Usually, the time window is divided into a few shorter logarithmically spaced gates (often 20), and the integral in eq. 4b is computed by sum of values achieved for each gate.

Apparent values are inverted for resistivity and chargeability using a Gauss-Newton iterative formulation (Loke and Barker 1996), where the chargeability dataset is inverted following the

6

144  linear approximation proposed by Oldenburg and Li (1994). We set inequality constraints on

145  chargeability ($M \geq 0.1$ mV/V) to avoid negative values in the inverted models. The goodness

146  of fit is evaluated for each line in terms of mean absolute percentage error for resistivity, while

147  for chargeability models is more convenient to use the mean absolute error, expressed in mV/V

148  (De Donno and Cardarelli, 2017).

149  The contribution of the surface conduction mechanisms, quantified by the normalized

150  chargeability $M_n$ (Slater and Lesmes 2002), can be calculated at the end of the inversion process

151  by:

152  $$M_n = \frac{M}{\rho} .$$  (5)

153

## 2.2  K-means clustering

155  K-means is a partitioning algorithm that divides the dataset into $K$ predefined non-overlapping

156  groups (clusters) of similar data. The similarity measure is based on a distance-based metric

157  that is typically the Euclidean distance. The algorithm starts assigning a set of $K$ means $\mu_1$,

158  $\mu_2,..., \mu_K$, also called centroids, computes the distance between each point and the $K$ centroids,

159  and assigns each point to the cluster $i$ with the closest centroid $\mu_i$. Actually, it assigns data

160  points to a cluster in order to minimize the sum of the squared distance (*SSE*) between the data

161  points and the cluster's centroids (Bhattacharya et al., 2021):

162  $$SSE = \sum_{k=1}^{K} \sum_{n=1}^{N} r_{nk} ||d_n - \mu_k||_2^2 ,$$  (6)

163  where $||..||_2$ is the Euclidean L2 norm, $d_n$ denote the data points, $N$ is the total number of points

164  in the dataset and $r_{nk}$ is a binary variable equal to is 1 if $d_n$ is assigned to cluster $k$ and 0

165  otherwise. Once each point has been assigned to a centroid, the procedure is iterated, i.e., the

166  centroid of each cluster is recomputed as the mean of all data points belonging to the same

167  cluster, and the category of each point is adjusted again until a maximum number of iterations

168  is reached, or the adjustment range is less than a given threshold.

7

169 K-means is a very powerful algorithm that, in addition to its simplicity, has the advantage of

170 having convergence guaranteed, as at each iteration step *SSE* always decreases (Bhattacharya

171 et al., 2021). However, K-means finds local minimum of *SSE* and different initial positions of

172 centroids determine different cluster solutions. To overcome this problem, the algorithm is run

173 many times placing the centroids in different random starting points, recording the variance at

174 each step and selecting the configuration corresponding to the minimum variance. Another

175 drawback of the algorithm is that the clustering depends on the number of *K* that needs to be

176 specified in advance. The optimal value of *K* is found by elbow method, if there is not any

177 geologic *a priori* information to constrain the number of clusters (Bhattacharya et al., 2021).

178 This method consists in varying the number of clusters *K*, then computing the percentage of

179 the explained variance *EV* for each trial:

180 $$EV_j = \frac{\sum_{i=1}^{j}(SSE_i - SSE_{i+1})}{\Delta SSE_{max}} \quad j = 1,2,\dots,K-1 \qquad (7)$$

181 where $\Delta SSE_{max}$ corresponds to the *SSE* deviation between 1 and the highest *K* among those

182 analysed.

183 The idea is that when increasing *K*, at some point the addition of another cluster does not

184 significantly improve the modelling of the data, and, thus, the best value of *K* is chosen as the

185 elbow of the *EV* curve (Thorndike, 1953). Our cluster analysis was performed by using

186 software packages available in the Statistics and Machine Learning Toolbox of MATLAB.

187

## 3. Field data

## 3.1 Case 1 (Southern Italy)

190     *3.1.1   Site location and geophysical measurements*

191 The survey area of Case 1 is located in the Campania region (southern Italy). Basically, no

192 information about the landfill design is available. This site is within a geological context

193 characterized by a dense alternation of layers of greyish silty clays and clay and of lithoid

8

194    arenaceous levels. This alternation is characteristic of the clayey-silty deposits largely

195    outcropping both in the landfill area itself and in the neighbouring areas (Urban Plan of

196    Montecorvino Pugliano, 2011).

197    Our geophysical survey consisted of resistivity and induced polarization measurements along

198    four profiles, each one 142.5 m long, oriented approximately SW-NE (Fig. 2), spaced about 50

199    m apart (Profile1 to Profile4). The data along each profile were collected simultaneously by

200    using a Syscal Pro 96 Switch resistivimeter (IRIS Instruments) in multi-electrode configuration

201    with a unit electrode spacing of 1.5 m. We used the pole-dipole array and, for each profile, we

202    positioned the (remote) current electrode 150 m away from electrode # 96, outside the landfill

203    area. Data were filtered for outliers, negative DC and/or IP voltage values or decay curves with
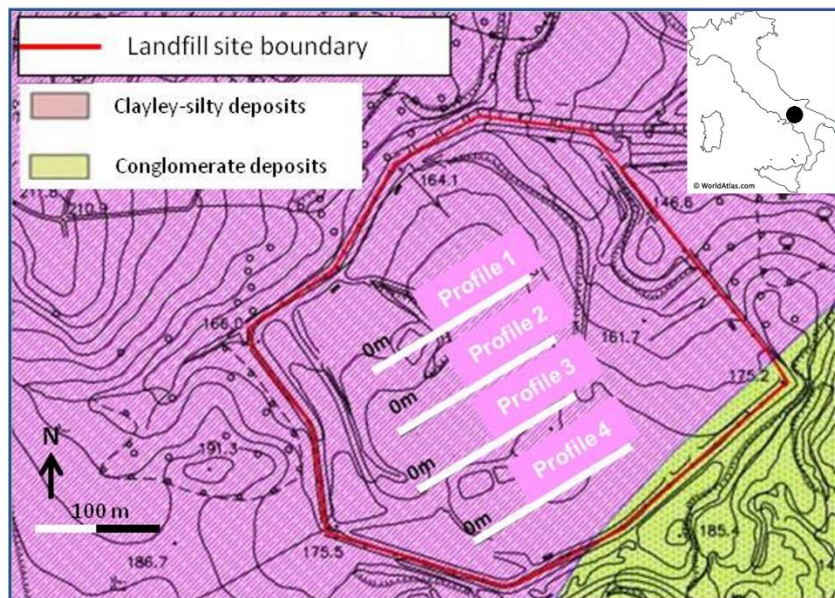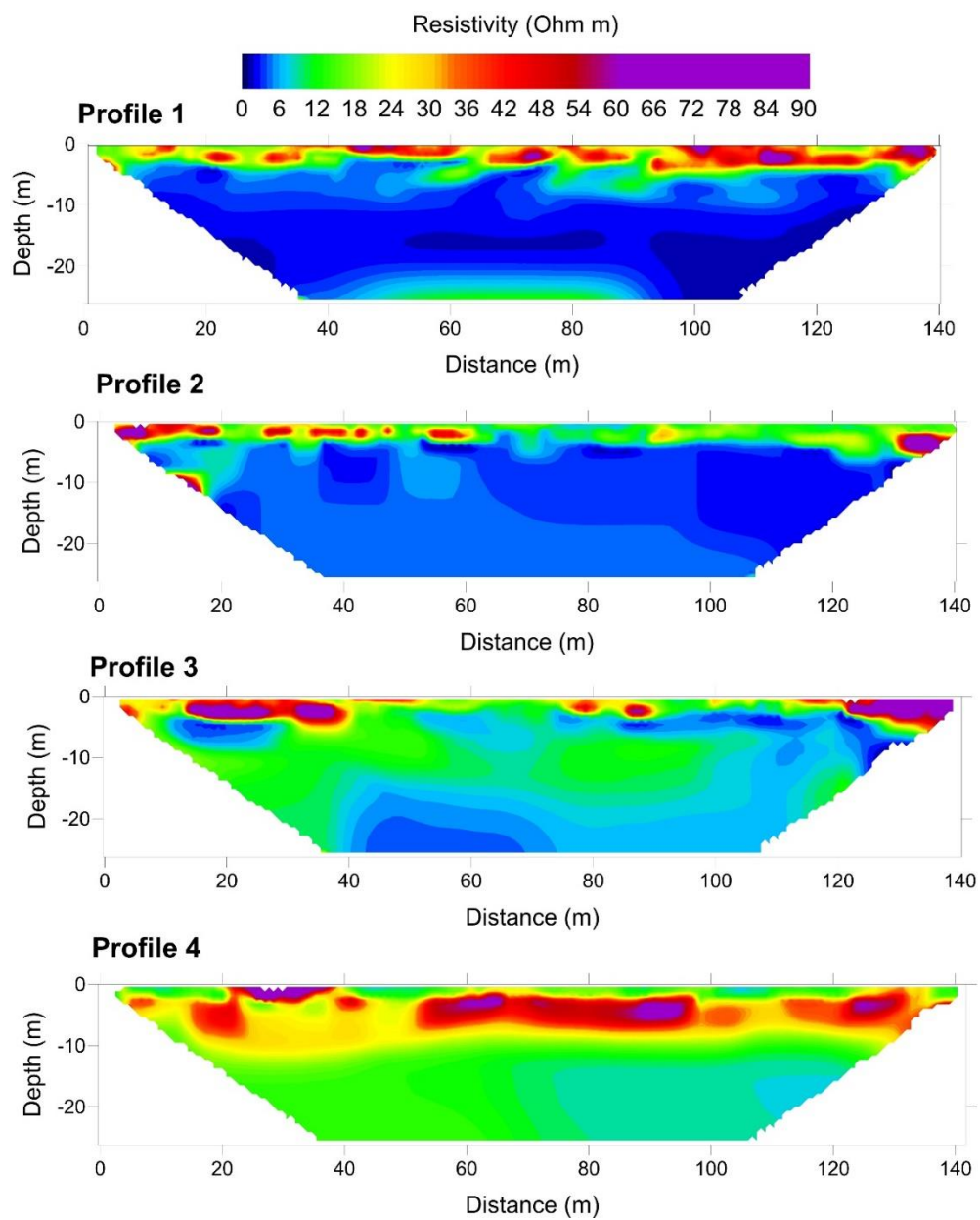
204    increasing voltage.



205

206    Figure 2. *Map of the survey area for the Case study 1 (landfill in Southern Italy). White lines*

207    *show the resistivity and induced polarization measurement profiles.*

208

209        *3.1.2   Inversion results*

210    The inverted models are shown in Fig. 3 (resistivity) and Fig. 4 (chargeability).

9

211

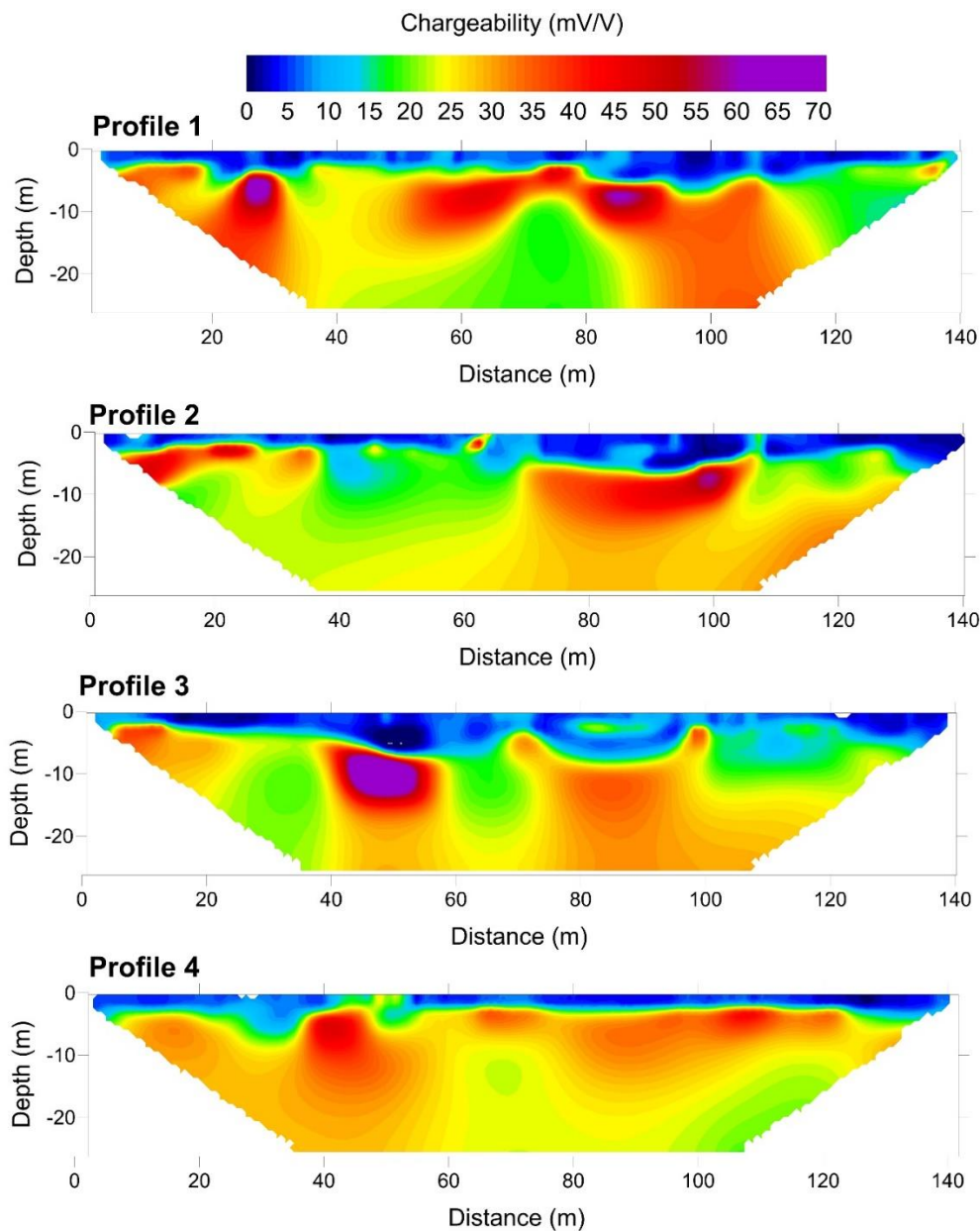Figure 3. *Resistivity models for the Case study 1*.

Figure 4. *Chargeability models for the Case study 1.*

For these four profiles we recognize a three-layer model from top to bottom:

- landfill covering and unsaturated waste, having a maximum thickness of 5 m, resistivity higher than 30 Ωm and chargeability close to zero, with local variation of thickness and resistivity (especially for Profile2 and Profile3) due to the high heterogeneity of the covering soil;

11

221      • saturated waste (leachate), with resistivity lower than 12 $\Omega$m and chargeability higher

222        than 20 mV/V. We notice strong lateral changes in resistivity and chargeability values

223        throughout the sections due to waste-changes in composition and degree of saturation.

224        The strongest geoelectrical anomalies ($\rho < 3$ $\Omega$m and $M > 35$ mV/V) are observed in

225        Profile1 that is the topographically lowermost profile, where an accumulation of

226        leachate is more likely. For the same reason Profile4 is the profile with the smaller

227        variations of resistivity and chargeability;

228      • a more resistive bottom layer, shown only in Profile1, whose resistivity (about 15 $\Omega$m)

229        is too low to be related to the presence of a bottom liner (geomembrane). The resistivity

230        and chargeability values of this bottom layer of Profile1 (values higher than 10 $\Omega$m and

231        chargeability of about 20 mV/V) suggest that it may be constituted by clayey-silty

232        deposits. For Profile1 we can thus estimate a maximum landfill thickness of about 20

233        m. The clayey-silty layer at the bottom of the landfill is very likely deeper for the

234        Profile2, Profile3 and Profile4 and therefore cannot be detected in our profiles.

235 Therefore (Figs. 3-4), we believe that leachate percolates downstream (from Profile4 to

236 Profile1) and accumulates at in the lowermost area, where Profile1 is located. Resistivity

237 models show a gradually decreasing value of resistivity from Profile1 to Profile4, whereas

238 chargeability models highlight for all profiles a high heterogeneity in leachate/waste

239 distribution.

240

## 3.2 Case 2 (Central Italy)

*3.2.1 Site location and geophysical measurements*

243 The investigated area of Case 2 is a landfill built in the '80s for being used as municipal waste

244 disposal of a medium-size city located in Central Italy. For this landfill, a few more information

245 about the original design is available. The landfill site (Fig. 5), located on a steep slope (slope

percent around 40%), was provided with a bottom liner (geo-membrane) overlying the in situ marly-arenaceous flysch. An embankment was built at the bottom (elevation around 445 m a.s.l.) to prevent slope instability phenomena. Geoelectrical investigation was planned for reconstructing the landfill depth and evaluating the leachate accumulation since it could reach a level that may trigger a slope failure. The supposed depth of the buried waste is greater upstream (southern part, elevation: 510-520 m a.s.l.), while reducing downstream (northern zone, elevation: 460-450 m a.s.l.).
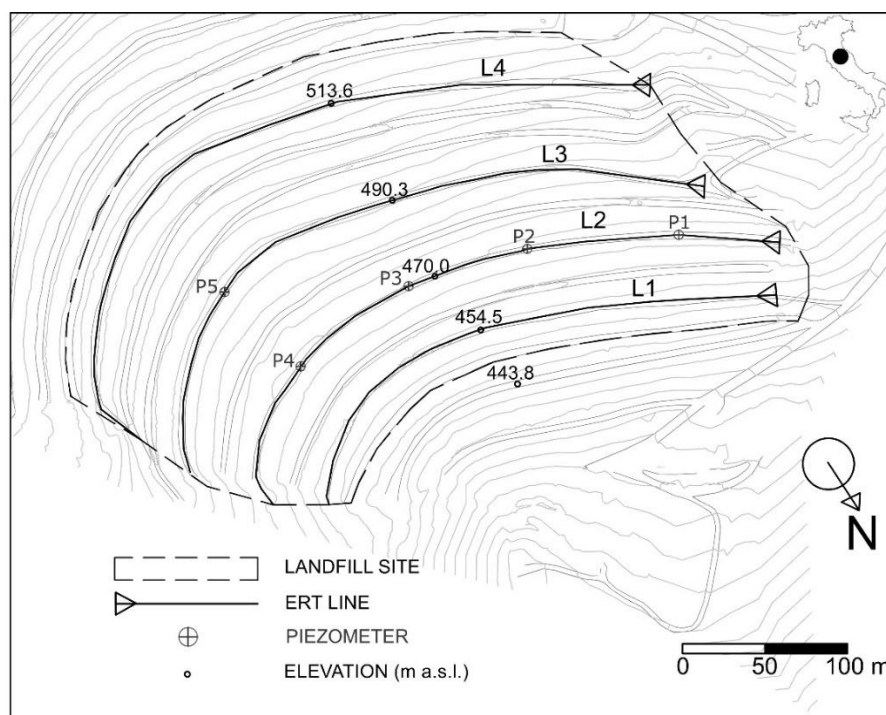


Figure 5. *Map of the survey area for the Case study 2 (landfill in Central Italy).*

The geophysical campaign encompasses four electrical profiles spaced approximately 40-50 m apart (L1 to L4), using the road tracks built for site management (Fig. 5). Five wells are located along L2 and L3 and piezometric levels were logged during the campaign to validate the geophysical models. Time-domain ERT and IP data were acquired by 48 electrodes spaced 5 m apart (Fig. 6), using the IRIS Instruments Syscal Pro resistivimeter. We employed the dipole-dipole array for data acquisition, using a maximum dipole length $a = 5$ and a maximum dipole

13

262   separation factor $n = 6$ (945 data points for each baseline), as it combines significant depth of

263   investigation and good lateral resolution needed to image the leachate variations. In fact, laying

264   cables outside the landfill for set-up a pole-dipole configuration was unfeasible in such a steep

265   and complex environment. For covering the long lines, we used the roll-along technique by

266   overlapping 36 electrodes for each baseline. We set the current injection time to 2 s (2 stacks),

267   a time delay of 40 ms and a logarithmic sampling of the IP decay curve using 20 gates (first

268   gate centered at 40 ms, last gate at 1.7 s). Similarly to what done for Case 1, data were filtered

269   for outliers, negative DC and/or IP voltage values or decay curves with increasing voltage.

270

271   *3.2.2   Inversion results*

272   The inverted models are shown in Fig. 6 (resistivity) and Fig. 7 (chargeability). The ERT and

273   IP sections were computed in 2.5D mode by linearization of the curvilinear profiles, to compare

274   the results with Case study 1. The error committed, calculated as the percentage deviation

275   between 2D and 3D geometric factors, is between 0.5% (L1) and 1.5% (L4), due to the low

276   curvature of profiles. We superposed the piezometric levels logged in the available wells to the

277   electrical models (white filled areas in Figs. 6 and 7), in order to validate the geophysical
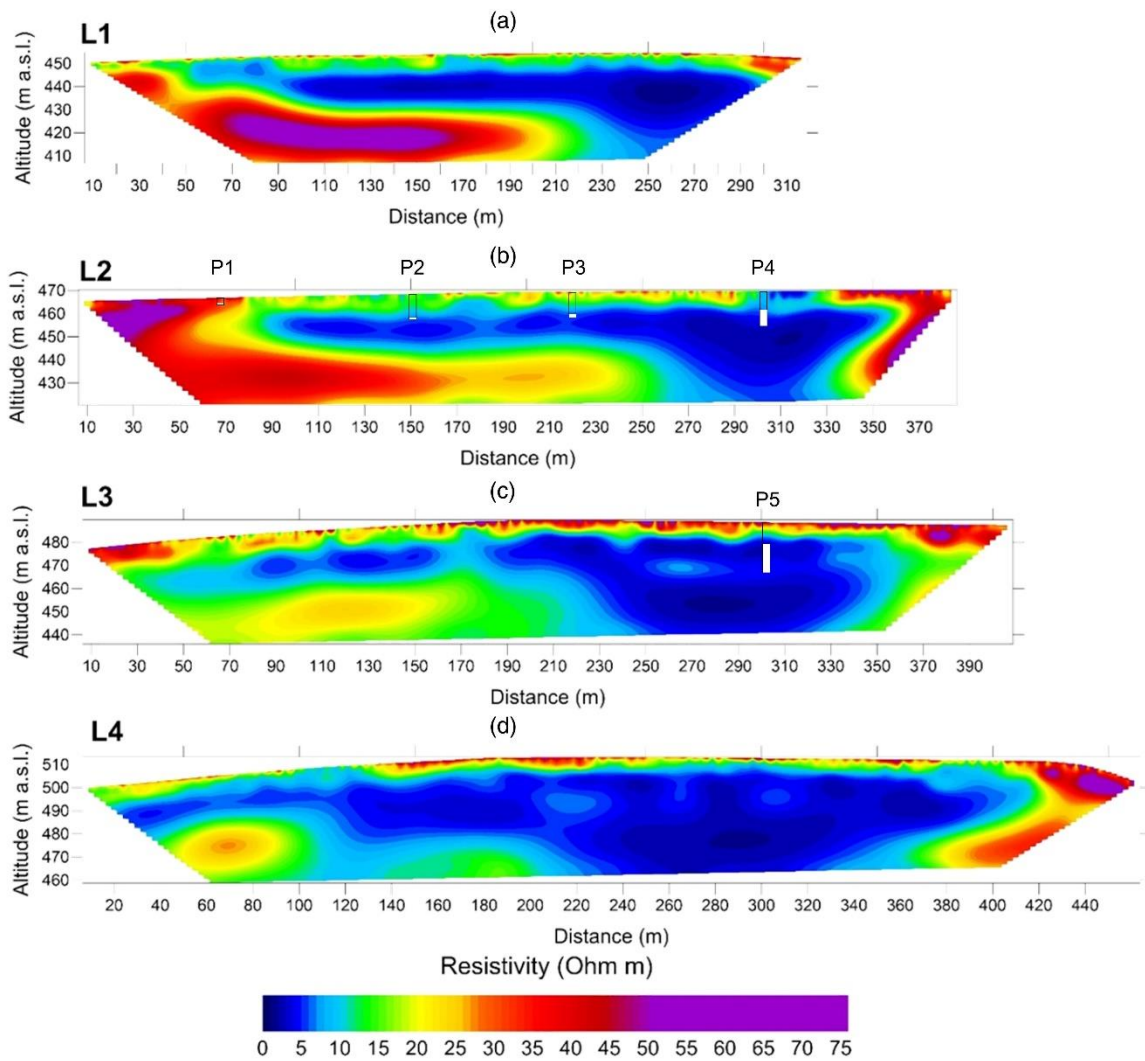
278   results.

279

Figure 6. *Resistivity models for the Case study 2: (a) L1, (b) L2, (c) L3, (d) L4. Absolute error is 2.2%, 4.2%, 7.4% and 4.5% respectively. The piezometric levels in wells (white filled areas) are superposed to the models.*
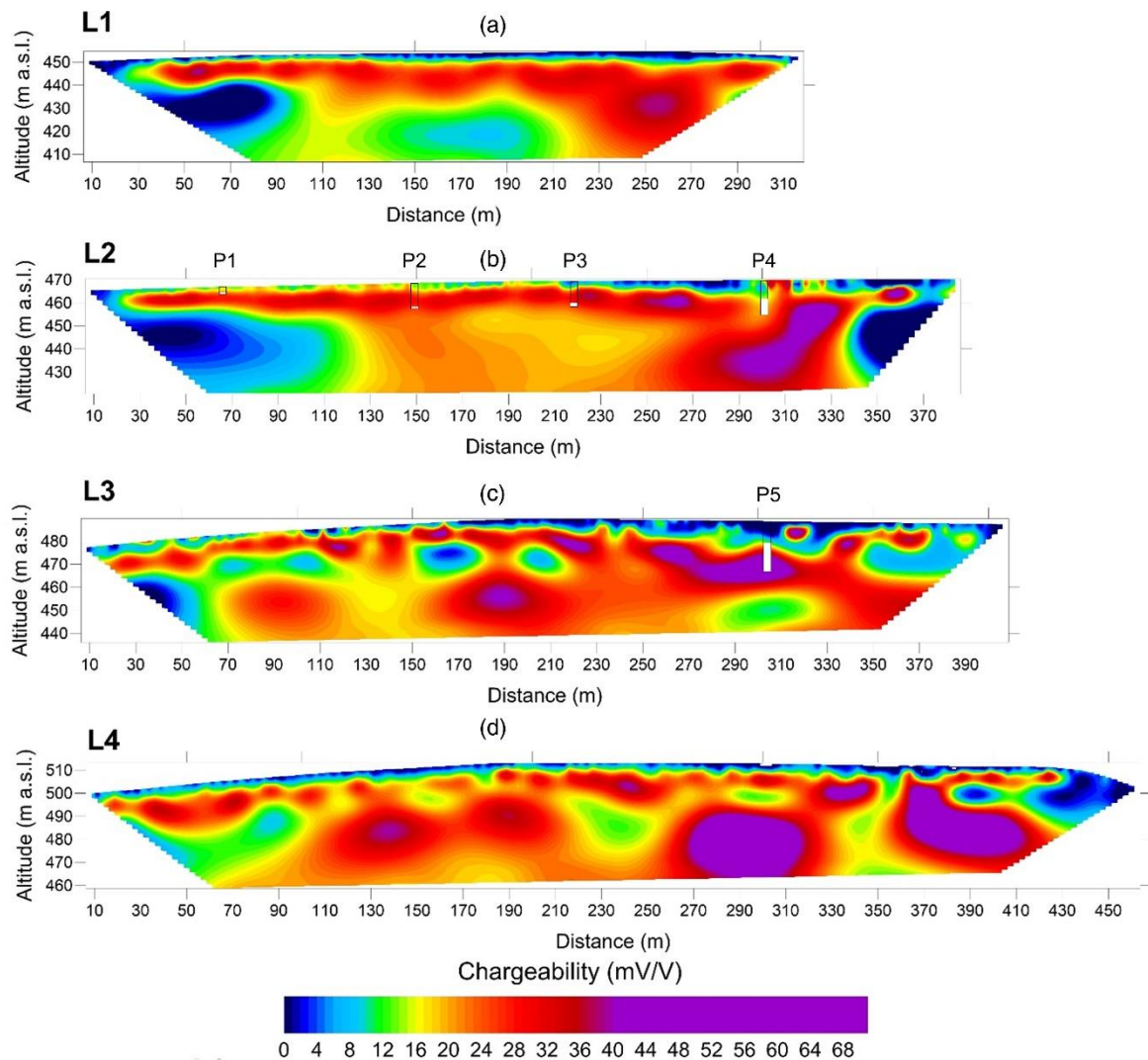
15

Figure 7. *Chargeability models for the Case study 2: (a) L1, (b) L2, (c) L3, (d) L4. Absolute*

*error is 1.1 mV/V, 2.1 mV/V, 3.7 mV/V, and 4.6 mV/V, respectively. The piezometric levels in*

*wells (white filled areas) are superposed to the models.*

Overall, we reconstruct for the four profiles a three-layer model from top to bottom:

- landfill covering and unsaturated waste, having an average thickness of about 8-10 m,

    resistivity higher than 15-20 $\Omega$m and chargeability close to zero, with local variation of

    thickness and resistivity due to the high heterogeneity of the covering soil;

- saturated waste (leachate), with resistivity lower than 10 $\Omega$m and chargeability higher

    than 10 mV/V. Strong changes in resistivity and chargeability values are clearly seen

16

296        throughout sections depending on composition and degree of saturation of the waste

297        mass. The strongest geoelectrical anomalies ($\rho < 3$ $\Omega$m and $M > 35$ mV/V) are located

298        in the deepest zones of the landfill between $x = 240 – 350$ m from the beginning of the

299        lines;

300      •   bottom liner (geomembrane) overlying bedrock (marly-arenaceous flysch), with

301        resistivity values higher than 10 $\Omega$m and chargeability lower than 10 mV/V. The

302        resistivity values are lower than expected for such a dielectric material (geomembrane),

303        because of the lack of sensitivity at greater depths (De Donno and Cardarelli 2017). A

304        series of sloped terraces is clearly visible, which ends in the deepest part of the landfill,

305        where we detect the strongest geophysical response. Consequently, the maximum depth

306        of the landfill can be estimated around 60-70 m upstream (L4) and 30-40 m downstream

307        (L1).

308 Therefore, leachate likely accumulates at the bottom of the landfill and percolates downstream

309 (from L4 to L1) through a preferential pathway, whose lateral extent progressively reduces

310 downstream from about 80 m (L4, between 240 to 320 m) to about 30 m (L1, between 240-

311 270 m). Chargeability models (Fig. 7) enhance the significant heterogeneity in leachate

312 accumulation, mainly for L3 and L4 lines (where accumulation zones are larger), because of

313 changes in degree of saturation or in waste composition. The piezometric levels in P4 and P5

314 wells confirms the geophysical reconstruction, while significant discrepancies can be noticed

315 between chargeability models and leachate levels for P1, P2 and P3. This residual ambiguity

316 leads the way for a more quantitative integration of geophysical data, which encompasses the

317 joint use of both inverted models as an input for the clustering analysis.

318

## 319   4. Results

## 320   4.1 Case 1 (Southern Italy)

321    The results of the K-means clustering analysis in the 3D parameter space defined by the

322    inverted values of $\rho$, $M$ and $M_n$ are shown in Figure 8. Since the inverted values span many

323    orders of magnitude, we applied a log10 transformation and then a normalization in the range

324    [0,1]. The clustering algorithm was iterated by varying $K$ from 1 to 50 and for each run; we

325    considered up to 500 different initial configurations of centroids choosing the configuration

326    that minimizes the distortion (i.e., the sum of point-to-centroid distances). We retrieved the

327    best choice for the number of clusters $K_{best}$=11 (Fig. 8a) using the elbow method with an

328    explained variance threshold equal to 95%.

329    The clustering analysis highlights 11 different regions in the parameter space, whose centroids

330    have the coordinates reported in Table 1 and are shown in Fig. 8b by blue markers. We

331    associated the cluster indices to a colour scale (green-yellow-red) by computing the Euclidean

332    distance to the point of normalized coordinates (0,1,1). This point corresponds simultaneously

333    to the lowest values of $\rho$ and the highest values of $M$ and $M_n$, and therefore can be associated

334    to the highest degree of contamination by leachate.

335    Being ERT and IP data acquired along the same profiles, triplets of values ($\rho$, $M$, $M_n$), in

336    addition to being characterized by a cluster index, are also associated to the same depths. Thus,

337    we use the retrieved cluster indices to get the integrated depth sections shown in Fig. 9. It is

338    worth noting that such cross sections do not represent a geoelectrical model in a traditional

339    sense, but combine information from all the investigated geophysical quantities.
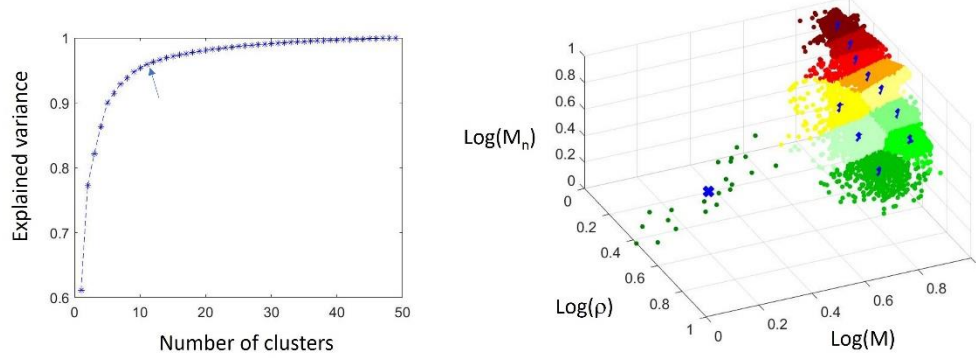
340

Figure 8. *Clustering analysis for the Case study 1: (a) Explained variance as a function of number of clusters. The best number of clusters is indicated by the arrow. (b) Scatterplot of geoelectrical data, where clusters are marked by different colours. In each cluster, the location of the centroid is shown by a blue cross.*

**Table 1.** Centroids coordinates of the clusters shown in Fig. 8b.

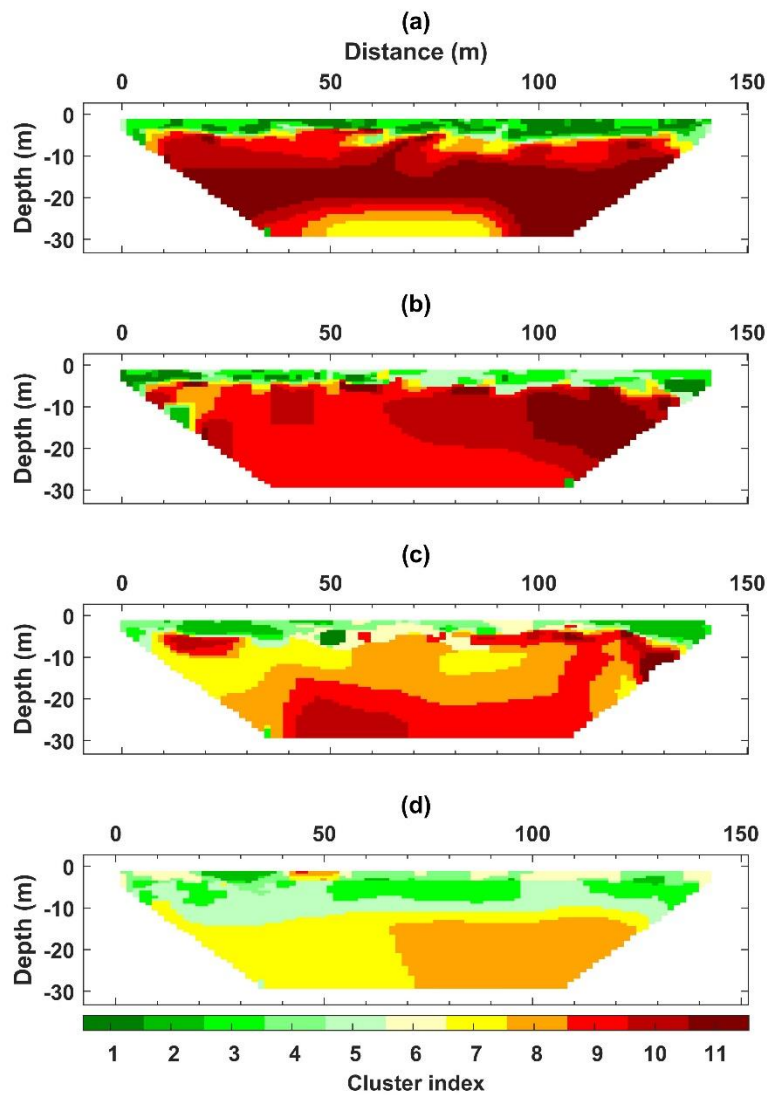| cluster color | $\rho$ ($\Omega$m) | $M$ (mV/V) | $M_n$ (mS/m) |
|---|---|---|---|
| dark green | 7.81 | 0.01 | 0.001 |
| | 46.54 | 4.13 | 0.09 |
| | 38.43 | 21.71 | 0.57 |
| ↓ | 19.42 | 4.62 | 0.24 |
| | 20.76 | 26.28 | 1.27 |
| light yellow | 9.14 | 4.92 | 0.54 |
| yellow | 11.23 | 26.74 | 2.38 |
| | 7.52 | 23.70 | 3.15 |
| ↓ | 4.82 | 23.21 | 4.82 |
| | 3.43 | 27.04 | 7.88 |
| dark red | 2.10 | 27.32 | 13.0 |

19

Figure 9. *Final cross-sections based on the proposed ML-based approach for the four profiles of the case study 1.*

## 4.2 Case 2 (Central Italy)

The results of the K-means clustering analysis performed on the second dataset (Figure 10) were obtained following the same clustering procedure as for Case 1. In this case, the 95% variance threshold determines the best configuration achieved with 10 different clusters (Fig. 10a). Looking at the distribution of data points in the parameter space (Fig. 10b), we note that the values of $M$ span over a wider range with respect to the Case study 1 (Fig. 8b). The darkest green cluster in Fig. 10b that groups data with a shape of straight line is related to the inequality

20

360  constraints set on the chargeability model during inversion to enforce positiveness of

361  chargeability (minimum value of 0.1 mV/V in this case).

362  The final output of the clustering algorithm are the cross-sections shown in Fig. 11, where we

363  observe a very good agreement between well and predictions of leachate contamination from
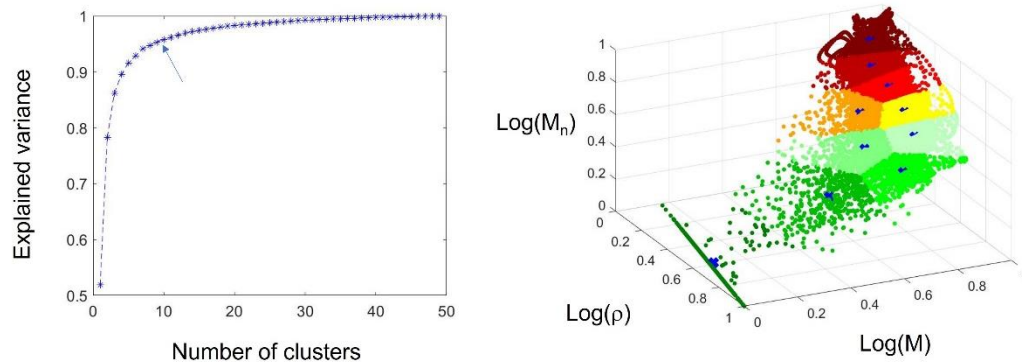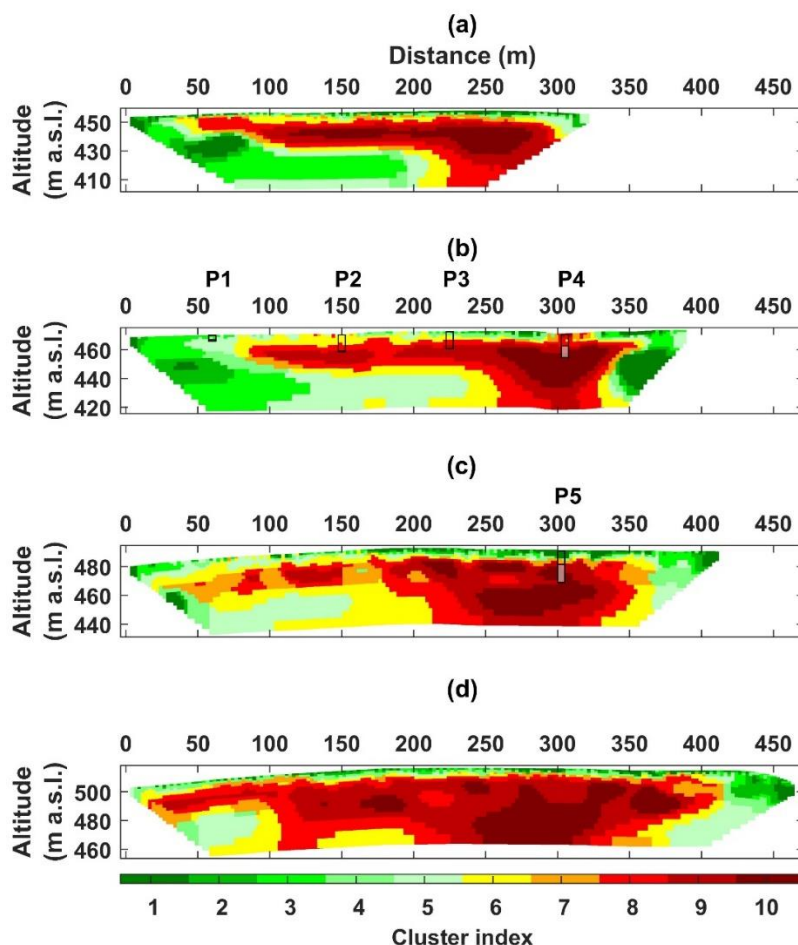
364  cluster analysis.



365

366  Figure 10. *Clustering analysis for the case study 2: (a) Explained variance as a function of*

367  *number of clusters. The best number of clusters is indicated by the arrow. (b) Scatterplot of*

368  *geoelectrical data, where clusters are marked by different colours. In each cluster, the location*

369  *of the centroid is shown by a blue cross.*

370

371  **Table 2.** Centroids coordinates of the clusters shown in Fig. 10b.

| cluster color | $\rho$ ($\Omega$m) | $M$ (mV/V) | $M_n$ (mS/m) |
|---|---|---|---|
| dark green | 31.77 | 0.11 | 0.003 |
|  | 30.30 | 1.89 | 0.06 |
| ↓ | 40.68 | 9.27 | 0.23 |
|  | 18.39 | 6.64 | 1.36 |
|  | 23.93 | 18.38 | 0.77 |
| yellow | 14.02 | 21.74 | 1.55 |
|  | 9.25 | 10.08 | 1.09 |
| ↓ | 7.77 | 23.56 | 3.03 |
|  | 4.59 | 23.19 | 5.05 |
| dark red | 2.83 | 32.37 | 11.44 |

21

372



(a)

Distance (m)

(b)

(c)

(d)

373

Figure 11. *Final cross-sections based on the proposed ML-based approach for the four profiles*

*of the Case study 2. The piezometric levels are superposed to the models.*

376

## 5. Discussion

To discuss and show the effectiveness of the proposed ML-based approach, in Figs. 12 and 13

we compare the $M_n$ models, achieved by the ratio of chargeability and resistivity of each pixel

in Figs. 4, 5 and 6, 7 (for the two study-cases respectively), with the sections retrieved by the

cluster analyses. In many cases the areas characterized by high values of $M_n$ fall within the

most hazardous zones identified by our cluster analyses (red clusters). Nevertheless, in many

other zones there are significant differences between $M_n$ models and the models obtained by

22

384    our ML-approach. It is worth noting the lack of lateral continuity of the chargeable zones (i.e.

385    Figs. 12a, 13a), compared to the respective ML images (Figs. 12e, 13e). This effect observed

386    in the $M_n$ sections increases the uncertainty on the model interpretation based only on inversion

387    results, preventing a clear detection of the leachate accumulation zones. The effectiveness of

388    K-means clustering in identifying different groups of data is particularly clear for the L2 and

389    L3 lines of Case study 2 (Fig. 13), where a quantitative comparison with data from wells is

390    available. The piezometric levels match well with the most hazardous areas (dark red clusters)

391    identified by our clustering analysis, whereas the levels are less correlated with the highest $M_n$

392    values. The very low level (0.2 m) logged in P1 is not highlighted by both procedures.

393    One of the main advantages of the proposed procedure is that the number of different zones is

394    directly retrieved from the cluster analysis and the shapes of such well-defined zones are not

395    affected by the choice of the colour scale. In fact, the scale was automatically distributed

396    according to the distance of the centroids from the point (0, 1, 1), which represents the

397    maximum level of contamination, given the electrical properties of leachate (highly conductive

398    and chargeable). We chose to scale colours from red to green as it is used in alert systems to

399    grade the severity of a hazard event. The final output is easy to interpret even for non-experts

400    in the field of geophysics. Typically, leachate is controlled by drilling monitoring wells and

401    our analyses show that the use of geophysical methods combined with ML techniques can

402    provide very accurate information on optimal locations to plan for such wells. The machine

403    learning-based approach can be therefore used to support decision-makers in the waste

404    management sector and reduce the costs of effectiveness of landfill management.

405    However, this study has some limitations and there is a room for further developments. In fact,

406    ML techniques are generally more effective if they use large datasets, and they can provide

407    different results if datasets are differently scaled. For this reason, before developing the

408    procedure, we preliminary explored different scenarios (Piegari and Paoletti, 2022; De Donno

409   and Piegari, 2022). We found that clustering gives better results if data are scaled and using a

410   3D parameter space instead of a 2D space with only $\rho$ and $M_n$. Additionally, the proposed

411   identification of the most contaminated zones is only based on the electrical properties of

412   leachate, since they are recognized as the most diagnostic for this purpose (Soupios et al.,

413   2017). However, the proposed multivariate analysis could be applied also to other geophysical

414   datasets, such as i.e. seismic tomography, which might add information about consolidation

415   and compaction of waste materials.



416

417   Figure 12. *Comparison between normalized chargeability sections (a-d) and integrated depth*

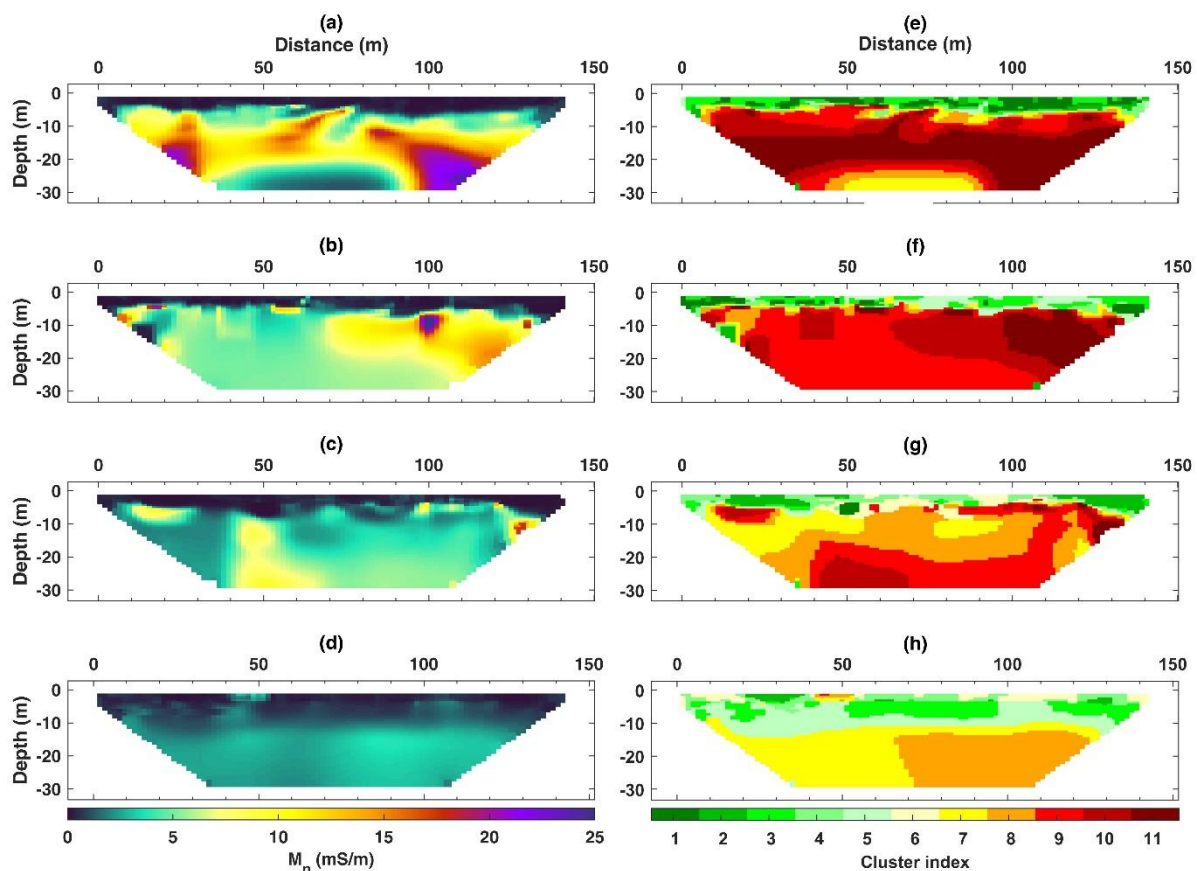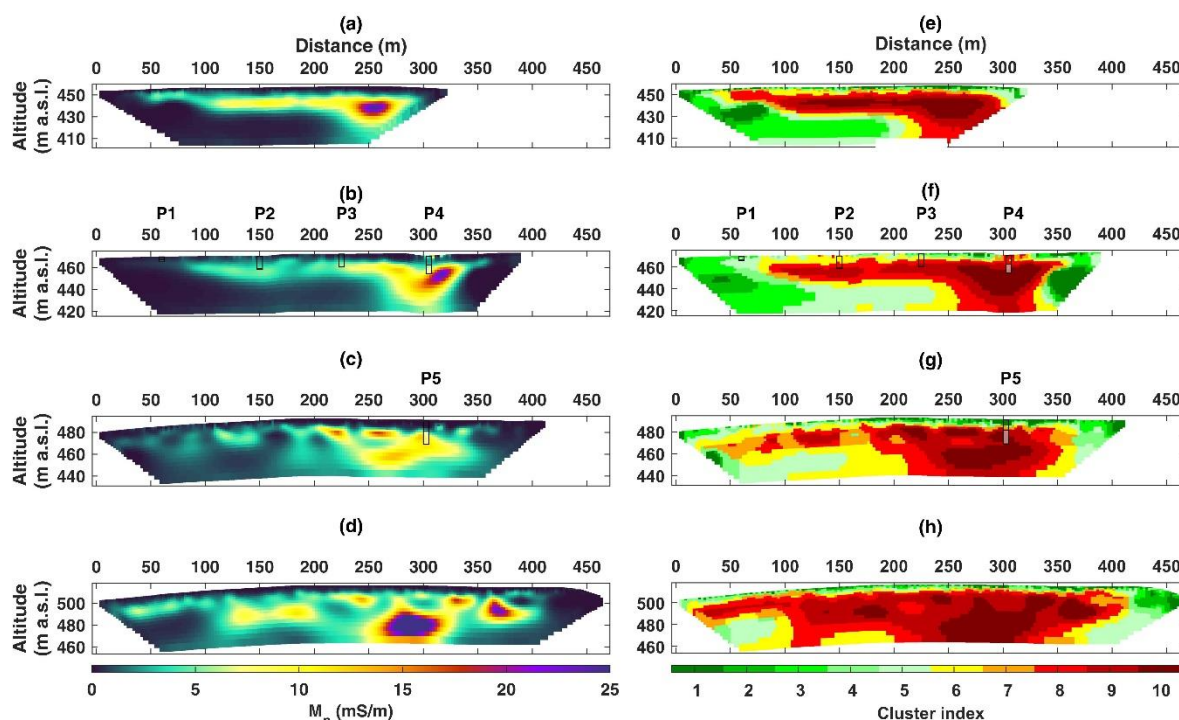418   *sections after clustering analysis (e-h) for the Case study 1.*

24

Figure 13. *Comparison between normalized chargeability sections (a-d) and integrated depth sections after clustering analysis (e-h) for the Case study 2.*

## 6. Concluding Remarks

In this study, we demonstrated that unsupervised machine learning may be successfully used to integrate data from resistivity and IP methods. The proposed ML-based approach is a flexible tool that can be easily adapted to other case studies also for different type of geophysical data. For each of the two investigated landfills, using K-means algorithm we were able to: i) achieve integrated model sections combining information from resistivity, chargeability and normalized chargeability data; ii) identify different regions in the investigated landfills associated with different leachate contamination levels; iii) locate the most hazardous zones. Therefore, our findings offer new perspectives for landfill characterization and have practical implications for landfill management. In fact, the final reconstructions of the investigated landfills help to predict the leachate flow pathways and enables more effective allocation of financial resources for the development of monitoring and remediation systems.

25

435

## Acknowledgments

441

## References

Abdideh, M., Ameri, A., 2020. Cluster Analysis of Petrophysical and Geological Parameters for Separating the Electrofacies of a Gas Carbonate Reservoir Sequence. Natural Resources Research 29(3), 1843–1856.

Bhattacharya, S., 2021. A Primer on Machine Learning in Subsurface Geosciences, 1st edn., Vol. 1, pp. 1–172, Springer.

Bernardetti, S., Bruno, P.P.G., 2019. The Hydrothermal System of Solfatara Crater (Campi Flegrei, Italy) Inferred from Machine Learning Algorithms. Frontiers in Earth Science 7, 286.

Binley, A., Slater, L., 2020. Resistivity and induced polarization: Theory and applications to the near-surface earth. Cambridge University Press.

Cesca, S., 2020. Seiscloud, a tool for density-based seismicity clustering and visualization. J Seismol 24, 443–457.

De Donno, G., Cardarelli, E., 2017. Tomographic inversion of time-domain resistivity and chargeability data for the investigation of landfills using a priori information. Waste Management 59, 302–315.

458    De Donno, G., Piegari, E., 2022. Clustering analysis of ERT/IP data for leachate mapping in

459       urban waste landfills. Near Surface Geoscience 2022, Belgrade, 18-22 September

460       (*accepted).*

461    Dey, A., Morrison, H.F., 1979. Resistivity modelling for arbitrarily shaped two-dimensional

462       structures. Geophysical Prospecting 27(1), 106-136.

463    Di Maio, R., Fais, S., Ligas, P., Piegari, E., Raga, R., Cossu, R., 2018. 3D geophysical imaging

464       for site-specific characterization plan of an old landfill. Waste Management 76, 629–642.

465    Ergene, D., Aksoy, A., Sanin, F.D., 2022. Comprehensive analysis and modelling of landfill

466       leachate. Waste Management, 145, 48-59.

467    Everett, M., 2013. Near-Surface Applied Geophysics. Cambridge: Cambridge University

468       Press. doi:10.1017/CBO9781139088435

469    Kamer, Y., Ouillon, G., Sornette, D., 2020. Fault network reconstruction using agglomerative

470       clustering: applications to southern Californian seismicity. Nat. Hazards Earth Syst. Sci. 20,

471       3611-3625.

472    Karpatne, A., Ebert-Uphoff, I., Ravela, S., Ali Babaie H., Kumar, V., 2019. Machine Learning

473       for the Geosciences: Challenges and Opportunities. IEEE Transactions on knowledge and

474       data engineering 31, 8, 1544.

475    Lavagnolo, M.C., 2019. Landfilling in developing countries. In: Cossu, R., Stegmann, R.

476       (Eds.), Solid Waste Landfilling: Concepts, Processes, Technologies, Elsevier, pp. 773–796.

477       https://doi.org/10.1016/B978-0-12-407721-8.00036-X.

478    Lindsey, C.R., Neupaneb, G., Spycher, N., Fairley, J.P., Dobson, P., Wood, T., McLing, T.,

479       Conrad, M., 2018. Cluster analysis as a tool for evaluating the exploration potential of

480       Known Geothermal Resource Areas. Geothermics 72, 358–370.

481    Lyra, G.B., Oliveira-Júnior, J.F., Zeri, M., 2014. Cluster analysis applied to the spatial and

482        temporal variability of monthly rainfall in Alagoas state, Northeast of Brazil. International

483        Journal of Climatology 34(13), 3546–3558, doi.org/10.1002/joc.3926

484    Loke, M.H., Barker, R.D., 1996. Rapid least-squares inversion of apparent resistivity

485        pseudosections by a quasi-Newton method1. Geophysical prospecting 44(1), 131–152.

486    Morita, A.K.M., Ibelli-Bianco, C., Anache, J.A.A., Coutinho, J.V., Pelinson, N.S., Nobrega,

487        J., Rosalem, L.M.P., Leite, C.M.C., Niviadonski, L.M., Manastella, C., Wendland, E., 2021.

488        Pollution threat to water and soil quality by dumpsites and non-sanitary landfills in Brazil:

489        A review. Waste Management 131, 163–176.

490    Mukherjee, S., Mukhopadhyay, S., Hashim, M.A., Gupta B.S., 2015. Contemporary

491        Environmental Issues of Landfill Leachate: Assessment and Remedies. Critical Reviews in

492        Environmental Science and Technology 45:5, 472-590.

493    Oldenburg, D.W., Li, Y., 1994. Inversion of induced polarization data. Geophysics 59(9),

494        1327-1341.

495    Piegari, E., Herrmann, M., Marzocchi, W., 2022. 3-D spatial cluster analysis of seismic

496        sequences through density-based algorithms. Geophysical Journal International 230, 2073-

497        2088, https://doi.org/10.1093/gji/ggac160

498    Piegari, E., Paoletti, V., 2022. Analysis of geoelectric data through machine learning

499        algorithms for waste leachate detection. Advances in Science, Technology & Innovation, *in*

500        *press*.

501    Power, C., Tsourlos, P., Ramasamy, M., Nirvolis, A., Mkandawire, M., 2018. Combined DC

502        resistivity and induced polarization (DC-IP) for mapping the internal composition of a mine

503        waste rock pile in Nova Scotia, Canada. Journal of Applied Geophysics 150, 40–51.

504    Raji, W.O., Adeoye, T.O., 2017. Geophysical mapping of contaminant leachate around a

505        reclaimed open dumpsite. Journal of King Saud University – Science 29, 348–359.

506     Shukla, S., Naganna, S., 2014. A Review on K-means DATA Clustering approach.

507        International Journal of Information & Computation Technology. 4(17), 1847–1860. ISSN

508        0974-2239

509     Seigel, H.O., 1959. Mathematical formulation and type curves for induced polarization:

510        Geophysics 24, 547–565.

511     Slater, L.D., Lesmes, D., 2002. IP interpretation in environmental investigations. Geophysics

512        67(1), 77–88.

513     Straus, D.M., 2019. Clustering Techniques in Climate Analysis. Climate Science, Oxford

514        https://doi.org/10.1093/acrefore/9780190228620.013.711

515     Soupios, P., Papadopoulos, N., Papadopoulos, I., Kouli, M., Vallianatos, F., Sarris, A., Manios,

516        T., 2007. Application of integrated methods in mapping waste disposal areas. Environ. Geol.

517        53(3), 661–675.

518     Soupios, P., Ntarlagiannis, D., Sengupta, D., Agrahari, S., 2017. Characterization and

519        monitoring of solid waste disposal sites using geophysical methods: current applications and

520        novel trends. Modelling Trends in Solid and Hazardous Waste Management, Edition 2017,

521        Chapter. fifth Ed. Springer, pp. 29. https://doi.org/10.1007/978-981-10-2410-8_5

522     Thorndike, R.L., 1953. Who belongs in the family? Pyschometrika 18 (4), 267–276.

523     Urban Plan of Montecorvino Pugliano, 2011. Official Bulletin of the Campania Region, No 1

524        of 3 January 2011. https://www.comune.montecorvinopugliano.sa.it/?page_id=788

525     Vaccari, M., Tudor, T., Vinti, G., 2019. Characteristics of leachate from landfills and dumpsites

526        in Asia, Africa and Latin America: an overview. Waste Manag. 95, 416–431.

527     WHO, 2015. Waste and human health: evidence and needs: WHO meeting report 5–6

528        November 2015: Bonn, Germany. https://apps.who.int/iris/handle/10665/354227

529    Zaini, M.S.I., Hasan, M., Zolkepli, M.F., 2022. Urban landfills investigation for leachate

530        assessment using electrical resistivity imaging in Johor, Malaysia, Environmental

531        Challenges 6, 100415.

532    Zhang, W., Zhang, Y., Gu, X., Wu, C., Han, L., 2022. Application of Soft Computing, Machine

533        Learning, Deep Learning and Optimizations in Geoengineering and Geoscience, 1st edn,

534        Vol. 1, pp. 1-138, Springer.