

1 **Cluster analysis of core measurements using heterogeneous data sources: an application to**
2 **complex Miocene reservoirs**

3

4 N. P. Szabó^{1,2*}, K. Nehéz^{1,3}, O. Hornyák^{1,3}, I. Piller^{1,3}, Cs. Deák¹, P. P. Hanzelik⁴, Cs. Kutasi⁴, K. Ott⁴

5 ¹University of Miskolc, 3515 Miskolc-Egyetemváros, Hungary

6 ²MTA-ME Geoengineering Research Group, University of Miskolc, 3515 Miskolc-Egyetemváros,
7 Hungary

8 ³HEIC, Higher Education Industrial Cooperation Centre, University of Miskolc, 3515 Miskolc-
9 Egyetemváros, Hungary

10 ⁴MOL Group, 1117 Október huszonharmadika Street 18, Budapest, Hungary

11 *corresponding author, Department of Geophysics, University of Miskolc, 3515 Miskolc-Egyetemváros,
12 Hungary, e-mail: norbert.szabo.phd@gmail.com

13

14 **Abstract**

15

16 An integrated clustering approach is suggested for the interpretation of petrophysical properties
17 measured on core samples. Porosity, carbonate content, bulk and grain density, permeability, irreducible
18 water saturation and capillary pressure curves measured from different sources are processed
19 simultaneously for a more reliable evaluation of hydrocarbon reservoirs. The specialty of the problem
20 is that the input dataset is composed of observations at different boreholes and depth intervals, where
21 the number of rock specimens and measurement types strongly vary. Several statistical methods such as
22 cluster and factor analyses are traditionally used for rock typing, which are either highly limited or not
23 feasible at all for the processing of such incomplete datasets. To overcome this difficulty, the sparse
24 matrix of multivariate observations is fully filled with reliable estimates of petrophysical parameters
25 before cluster analysis. In the first phase, a correlation-based interpolation method is proposed for
26 replacing the missing data with synthetic ones estimated from the available petrophysical information.
27 Then, principal component analysis of the filled data matrix is performed to investigate the relative
28 contribution of different variables to the solution and reduce the large number of measured parameters

29 into fewer variables. In the last step, a non-hierarchical cluster analysis of the principal components is
30 made to separate the lithological units and reservoir zones. The suggested statistical workflow is tested
31 in a Hungarian oilfield, where Miocene reservoirs of different lithologies are evaluated using a large
32 amount of laboratory data collected for three decades. When clustering all petrophysical variables in a
33 joint procedure, not just the lithological properties but also the fluid and other reservoir characteristics
34 are taken into account to differentiate the pay zones from unproductive intervals. In addition to the
35 current application, the statistical method may serve as useful tool for improved well log analysis, well-
36 to-well correlation and reservoir modeling on larger scales.

37

38 **Keywords** K-means clustering; petrophysical properties; sparse matrix; core measurement; Miocene
39 reservoir; Hungarian oilfield

40

41 **1. Introduction**

42

43 The role of multivariate statistical approaches has been continuously increasing with the growing
44 amount of observations and highly developing computational resources in petroleum geosciences.
45 Hempkins (1978) suggests the practical use of multiple regression, cluster and principal component
46 analysis as powerful tools in formation evaluation. For today, these techniques became as routinely
47 applied methods in petrophysics, used primarily for the identification of hydrocarbon reservoirs, facies
48 analysis, determination of rock physical relations, trends and correlation analyses, and for the
49 replacement of missing data along the boreholes. Principal component analysis and its extended variant
50 called factor analysis were originally developed to reduce the dimension of the data space defined by
51 the measurement variables (Lawley and Maxwell, 1962), the geological application of which makes it
52 possible to emphasize the common geological/geophysical information in measured dataset and explore
53 hidden variables dependent on lithological and rock physical characteristics that cannot be measured
54 directly with geophysical instruments. A modern application of factor analysis can be found in Jarzyna
55 et al. (2017), where the heterogeneity of Paleozoic shale gas formations and the complex geophysical
56 responses of organic-rich reservoirs were studied. The robust forms of the above exploratory statistical

57 methods make a significant improvement in the quality of parameter estimation. Szabó and Dobróka
58 (2017) processed oilfield well logs by iteratively reweighted factor analysis to give an outlier-free
59 estimate to the shale volume of clastic formations. Permeability as a related quantity was also predicted
60 by the same methodology using an evolutionary computation-based optimization approach (Szabó and
61 Dobróka, 2018).

62

63 Cluster analysis being an effective rock typing tool has a wide literature record. As an advanced
64 example, Skalinsky et al. (2006) suggests the clustering of mercury injection and well-logging data to
65 predict the types of carbonate rocks. Beyond its capabilities with stratigraphic evaluation of rock
66 formations, cluster analysis can provide the inversion of well logging data with quantitative information
67 such as matrix and fluid parameters of the separated lithological units (Szabó et al., 2013). A clustering-
68 based uncertainty assessment of porosity estimation is proposed by Masoudi et al. (2018), where the
69 accuracy information of the former parameter can be extended to predict irreducible water saturation
70 and permeability by means of fuzzy arithmetic. In the petrophysical characterization of unconventional
71 reservoirs, cluster analysis seems to be also an inescapable tool (Ma and Holditch, 2016). Moreover,
72 soft computing methods including self-organizing maps, neural networks, machine learning techniques
73 and other artificial intelligence approaches will be more and more important in the near future of
74 hydrocarbon exploration (Cranganu et al., 2015).

75

76 However, the performance of clustering spectacularly decreases in big and incomplete datasets.
77 Classical non-hierarchical clustering algorithms are known to be very sensitive to the initial setting of
78 cluster centers, which may cause interpretation problems in the lack of a priori geological information.
79 On the other hand, traditional methods would either drop the data columns of the missing data, or fill
80 the missing data with zero values. The former approach dramatically reduces the amount of data
81 involved in clustering, the latter may introduce false data. Thus, neither of the above methods provide
82 feasible solution in petrophysical applications. Considering the aforementioned level of incompleteness,
83 it is expedient to estimate missing data before applying data mining techniques.

84

85 Several matrix completion and imputation algorithms exist for handling missing data. One of the
86 simplest approaches replaces the missing entries with the mean (or median) of each column of the input
87 data matrix (Little, 1986). As another alternative, the nearest neighbour algorithm weights the sample
88 elements using the mean squared difference of columns for those two rows where there are observed
89 data. In our application, there are huge blocks of missing sub-matrices because of the separated groups
90 of unknowns that are available in different wells. In this case, the traditional methods do not work
91 properly. Recently, newer soft-impute, nuclear norm minimization and classical direct matrix
92 factorization algorithms have been introduced to carry out the completion of missing values. The soft-
93 impute method which iteratively replaces the missing elements with those obtained from a soft-
94 thresholded singular value decomposition did not yield satisfactory results for our problem (Mazumder
95 et al., 2010). Nuclear norm minimization offers an exact matrix completion via convex optimization
96 (Candès et al., 2009). The direct matrix factorization method transforms the incomplete matrix into an
97 M -by- N matrix with an L_1 sparsity penalty on the elements of the M -th row and an L_2 penalty on the
98 elements of the N -th column using gradient descent algorithm (Gemulla et al., 2011). By using the above
99 methods, serious problems were raised in the physical interpretation of completed data matrices. Our
100 preliminary investigations showed the application of the above mentioned interpolation algorithms
101 results in a matrix with wrong data, for example negative or outlying values were frequently obtained
102 being out of the expected physical ranges of petrophysical variables.

103

104 In this study, an efficient approach is proposed to eliminate the disadvantages of the clustering
105 procedure. By integrating numerous variables measured on core samples in the laboratory, which are
106 based on different physical principles, we aim to increase the performance of clustering and resolve the
107 problem of ambiguity. When doing this, we transform the sparse matrix of measured variables into a
108 fully filled one by replacing the missing data by quasi measured data using a novel regression-based
109 (multivariate) imputation method. Cluster analysis of the data matrix completed with the synthetic
110 values of petrophysical quantities gives highly reliable results for rock typing and allows an improved
111 identification of reservoir zones, which is demonstrated by a Hungarian oilfield study.

112

113 2. Methods

114

115 2.1 Imputation algorithm

116

117 To perform a reliable imputation of core laboratory data, a correlation-based algorithm is proposed in
118 the paper, which effectively predicts the missing values from the available petrophysical observations.

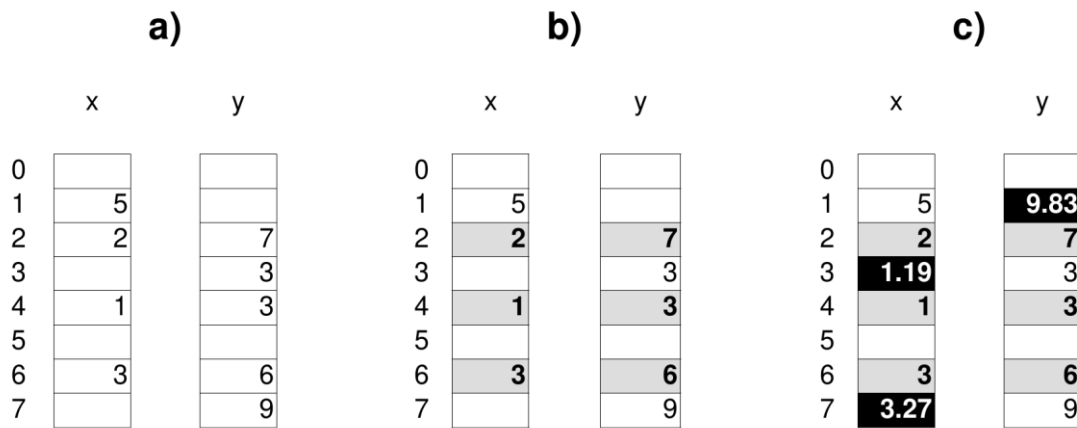
119 The heterogeneity of data sources is a result of the measurements taken at different times, during which
120 sometimes new measurement types are added or older ones are taken away due to the development of
121 measuring technology. The change of the observed parameters with time may also be the result of the
122 re-interpretation of geological formations, which requires modification of the measurement program. It
123 is expedient to maximize the number of measured properties (i.e., the number of the data columns),
124 because the larger the statistical sample the more reliable the result is, revealing more characteristics of
125 the studied formations.

126

127 The correlation-based imputation method assumes a sparse data matrix as input, where places of missing
128 data are filled with NaN (Not-a-Number) values. It is noted that the replacement of all NaN values is
129 not always possible. The proposed algorithm shows a viable solution for the datasets used in practice,
130 but the limitations of the described heuristic approach require further investigation. It is assumed that
131 the columns of the matrices have been sorted by significance, which is not an exact measure in this case.
132 The assumption is based on the fact that the data matrices have been arranged by petrophysical experts.
133 They tend to choose lower column indices for more important properties, while the higher indices remain
134 for the least important ones. The imputation algorithm uses this implicit information by filling the
135 missing data from lower indices to higher indices in multiple steps. As a basic calculation step of the
136 estimation process, one fits a linear regression model for two selected columns. Consider variables x
137 and y as two columns of the input data matrix. First, we collect (x_k, y_k) pairs of available values (where
138 x_k and $y_k \in \mathbb{R}$). It is important to check whether the number of points is enough for linear regression.
139 (The examined core datasets are proper in the sense that we have enough points for linear regression in
140 all cases.) In positive case, we can fit a linear regression model using those rows where we have data in

141 both columns. Then, one can give an estimate for the missing values by substituting the data into the
 142 linear equation, where exactly one of the x_k or the y_k is missing (NaN). The imputation process is shown
 143 in Fig.1. The selected columns of variables x and y with four and five measured data, respectively, are
 144 represented in Fig.1a. The white cells denote the NaN values, where measurements are not available. In
 145 Fig. 1b, the (x_k, y_k) pairs with real values are highlighted. This example has only three matching points
 146 for linear regression, but the number of data points is significantly larger in the processed big dataset.
 147 Regression analysis finds an optimal linear function, the parameters of which are estimated by
 148 minimizing the sum of squared distances of the selected points from the line. The above described
 149 algorithm is applied to each row of the input data matrix until it is fully filled with values extracted from
 150 the resultant linear equations (Fig. 1c). The pseudo-code of the correlation-based interpolation algorithm
 151 is presented in the Appendix.

152



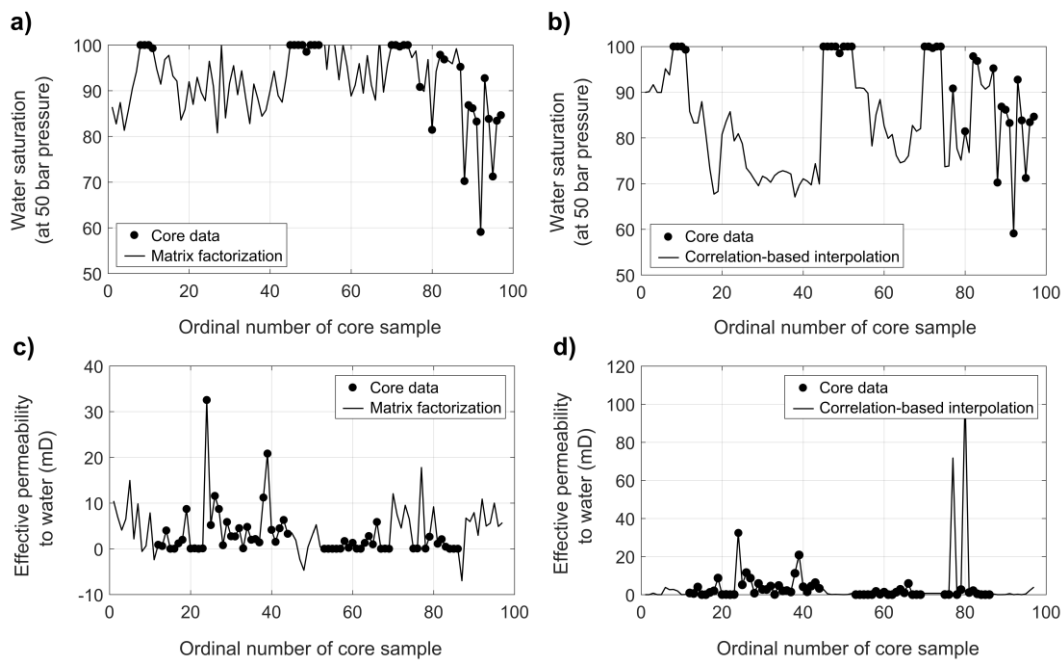
153

154

155 **Fig. 1.** Scheme of the correlation-based imputation method as a data preparation and filling procedure
 156 before cluster analysis
 157

158 A comparison can be made between the result of matrix factorization as a traditional completion method
 159 (see Section 1) and the proposed imputation procedure for a Miocene dataset (Fig. 2). On the axis of
 160 abscissae, the number of rows of the input data matrix is plotted representing the depth of core samples
 161 collected from neighboring wells. As it is seen in Fig. 2, the matrix factorization algorithm, on the
 162 examined data, overestimates the water saturation (~ 0.9 v/v) in high porosity gas-bearing formations. It
 163 estimates too low gas saturation or forecasts water-bearing formations furthermore it gives noisier

164 solution than the correlation-based imputation algorithm (Fig. 2a–b). In addition to it, non-physical
 165 estimates for the petrophysical parameters are also given, e.g., negative values of permeability are
 166 indicated in Fig. 2c. The reliability of matrix factorization is questionable, because it overestimates the
 167 permeability of calcareous marls and silts (e.g., for the first eight core samples in Fig. 2c) and at the
 168 same time, it underestimates the same property in fractured basalt formations (compare the 77th and 80th
 169 core samples in Fig. 2c–d). The permeability values in fractured zones do not differ significantly from
 170 their environment by means of classical matrix factorization, while they actually do separate using our
 171 proposed method (see the two peaks above 60 mD in Fig. 2d). In contrast, the correlation-based
 172 imputation method gives more realistic results and physically correct values of petrophysical
 173 parameters.
 174



175
 176 **Fig. 2.** Interpolation results for the capillary pressure derived water saturation (a–b) and effective
 177 permeability (c–d) using a traditional matrix factorization method (a, c) and the suggested correlation-
 178 based imputation procedure (b, d) in a Hungarian Miocene formation
 179

180 2.2 Dimension reduction

181
 182 In order to make the interpretation of multivariate measurements easier, it is advantageous to reduce the
 183 number of observed parameters to fewer statistical variables (Jolliffe, 2002). In this study, principal

184 component analysis (PCA) is used to decrease the number of observed petrophysical variables by
 185 reducing the number of columns of the input data matrix into smaller sizes. In complex oilfield problems,
 186 PCA can be beneficially used as a preliminary data processing step to find the main geological
 187 characteristics of the reservoir model. As an example, Jung et al. (2018) established a PCA assisted
 188 support vector machine approach for setting a proper initial model, which allows an improved history
 189 matching and prediction of reservoir performances in heterogeneous channel reservoirs. The PCA
 190 transformation is orthogonal, which gives new uncorrelated variables, i.e., principal components (PCs).
 191 Consider \mathbf{D} as an N -by- M fully filled data matrix, where N denotes the number of core samples collected
 192 in different wells and M is the total number of observed petrophysical variables in all wells. Let us write
 193 the data matrix as a product of the N -by- r matrix of PCs (\mathbf{P}) and that of their coefficients of M -by- r size
 194 (\mathbf{W})

195

$$196 \quad \mathbf{D} = \mathbf{P}\mathbf{W}^T \quad (1)$$

197 where r is the number of PCs. With the knowledge of the PCs' coefficients, a unique solution can always
 198 be found to Eq. (1). Scores of the j -th principal component are given by elements P_{lj} , where $l=1,2,\dots,N$
 199 and $j=1,2,\dots,r$. By multiplying Eq. (1) from the right with an orthonormal matrix \mathbf{W}^T , the PCs can be
 200 determined as a linear combination of the measured variables by equation $\mathbf{P}=\mathbf{D}\mathbf{W}^T$. The estimation of
 201 the weighting coefficients leads to the solution of an eigenvalue problem, in which the covariance matrix
 202 of the matrix product $\mathbf{D}^T\mathbf{D}$ for centralized data gives

203

$$204 \quad \lambda_k = \sigma_k^2, \quad (2)$$

205

206 where λ_k is the k -th eigenvalue and σ_k is the variance of the k -th petrophysical parameter. As a
 207 consequence of Eq. (2), the extracted PCs are sorted in a manner that the first few of them explaining
 208 most of the variance of the original variables showing the directions of biggest variances of the original
 209 sample. PCA specifies the coordinates of the objects of observed data in a new coordinate system
 210 stretched by the PCs, and rotates the original variables into the direction of the principal axes. The

211 number of PCs are arbitrarily chosen or selected by studying the distribution of eigenvalues in the
212 function of the PCs. By examining the elements of the M -by- M loading matrix \mathbf{W}^T , one can investigate
213 the individual contribution of measurement variables to the PCs. In this study, the relative importance
214 of the i -th petrophysical variable is expressed by

$$216 \quad s_i = \max_j |W_{ij}|, \quad (3)$$

217
218 where s gives the vector of maximal absolute values of PCs' weights ($i=1, 2, \dots, M$ and $j=1, 2, \dots, r$). In
219 the reduced coordinate system of PCs, cluster analysis can be made to group the data objects and infer
220 petrophysical characteristics of the available core information.

221

222 2.3 Clustering of big core datasets

223

224 A non-hierarchical K-means clustering was applied to the interpolated core data matrix, which is a
225 commonly used clustering approach for performing unsupervised learning tasks (Hartigan, 1979). It can
226 be effectively used for the grouping of core data in such a way that the M -dimensional objects specified
227 by petrophysical properties measured on given rock samples (collected from given depths and wells)
228 are more similar than others observed on different samples. From the point of view of the proposed
229 method, it is of great importance that data objects connected to the same cluster define approximately
230 the same lithological and petrophysical character, while other clusters represent dissimilar ones. Given
231 an initial set of center groups calculated as the average of cluster elements, the algorithm proceeds by
232 repeating two steps. We assign each observation to the cluster by the nearest mean principle. Then, we
233 calculate the new means to be the centroids of the observations in the new clusters. The solution is
234 obtained when assignments no longer change. The optimal number of clusters (K) to be formed can be
235 selected by the minimization of SSE (Sum of Squared Error), which gives the deviation between the
236 centroid and group elements for all groups. In the preliminary determination of the number of clusters,
237 it is important to consider a priori geological and rock physical information. Too small number of

238 clusters does not allow proper spatial resolution, but a large number of groups may result in false
 239 conclusions, i.e., non-existent lithological categories.

240

241 Another factor influencing the result of cluster analysis is the distance between the data objects as the
 242 degree of similarity. Let $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ two vectors having M number of observed variables in the data space
 243 of X_1, \dots, X_M . In our case, X_p denotes the p -th petrophysical variable measured in the laboratory. In a more
 244 detailed form, the i -th and j -th objects are given as $\mathbf{x}^{(i)} = [x_1^{(i)}, \dots, x_M^{(i)}]^T$ and $\mathbf{x}^{(j)} = [x_1^{(j)}, \dots, x_M^{(j)}]^T$ (where T
 245 is the symbol of transpose). The most commonly used distance metric is the Euclidean norm

246

$$247 \quad D_E(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \left[(\mathbf{x}^{(i)} - \mathbf{x}^{(j)})^T (\mathbf{x}^{(i)} - \mathbf{x}^{(j)}) \right]^{1/2}, \quad (4)$$

248

249 which may be very sensitive to outliers that are extremely far from the center of the group. A more
 250 robust solution can be achieved by using other distance metrics such as the L₁-norm based City block
 251 (or Manhattan) distance

252

$$253 \quad D_C(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \sum_{p=1}^M |x_p^{(i)} - x_p^{(j)}|, \quad (5)$$

254

255 or the Mahalanobis distance taking the correlation of the data into account

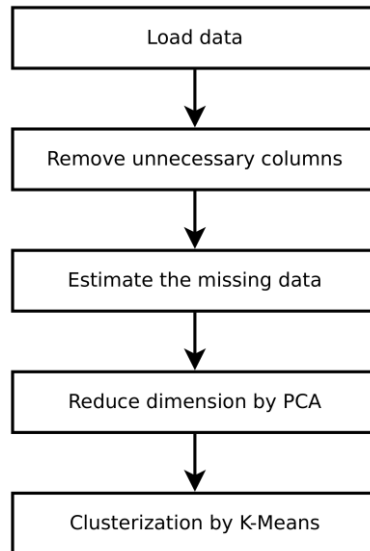
256

$$257 \quad D_M(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \left[(\mathbf{x}^{(i)} - \mathbf{x}^{(j)})^T \mathbf{S}^{-1} (\mathbf{x}^{(i)} - \mathbf{x}^{(j)}) \right]^{1/2}, \quad (6)$$

258

259 where \mathbf{S}^{-1} is the inverse of sample covariance matrix being a weighting operator during the cluster
 260 analysis of data objects. By combining the process of K-means cluster analysis with PCA, the following
 261 simplified workflow shown in Fig. 3 was developed.

262



263

264 **Fig. 3.** Flowchart of the correlation-based interpolation- and principal component analysis assisted
 265 clustering method
 266

267 In applying the method, we have considered a column as unusable, when it contains only the same
 268 values. One can drop these kind of columns without information loss. It must be also mentioned that
 269 clustering can be applied with or without dimension reduction. We prefer the use of PCA method in
 270 most cases, because it makes possible to highlight the most important features. Unfortunately, it does
 271 not provide the meaning of the resulted new dimensions. However, the physical relations between the
 272 PCs and petrophysical characteristics may be explored by partial correlation analysis, which can be
 273 added as a new element to the workflow.

274

275 3. Case study

276

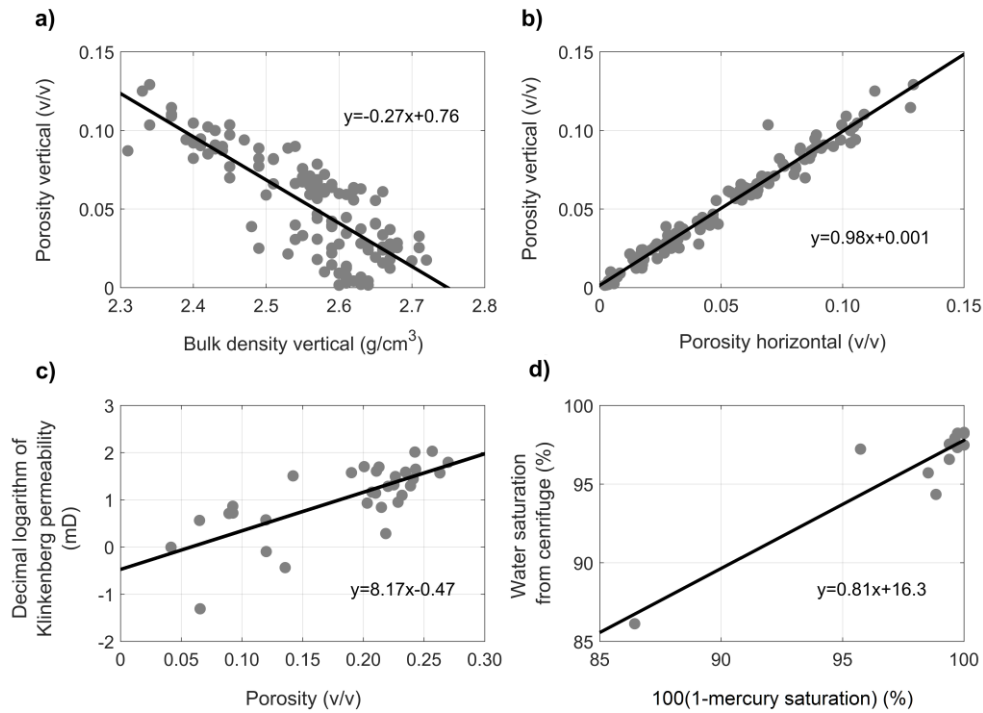
277 The core samples examined in this study originate from the southern part of the Pannonian Basin
 278 Province of Central Europe, where several petroleum systems have been discovered. The Pannonian
 279 Basin consists of a large extensional basin of Neogene age overlying Paleogene basins and a Mesozoic
 280 (or older) basement (Dolton, 2006). A several kilometers thick large-extension Tertiary basin-fill
 281 sedimentary sequence contains oil and gas-bearing formations. The main lithological categories are
 282 conglomerate, breccia, calcareous marl, clay, aleurolit, and gravellous sandstone of different grain sizes,
 283 dolomite and some basalt. The rock specimen was collected in 28 neighboring wells from an interval of

284 1352 m. The full dataset as input for cluster analysis includes 421 core samples and 59 measured
 285 petrophysical properties, i.e. carbonate content, bulk and grain density, helium porosity, Klinkenberg-
 286 corrected permeability, oil and water permeability, irreducible water saturation, sample porosity and
 287 volume, mercury saturation in the pressure range of 0.1–4,000 bar (and pore-throat radii between
 288 $1.88 \cdot 10^{-3}$ –75 μm), and water saturation determined by centrifuge method in the pressure range of 0.16
 289 and 6 bar. Most of the observed variables were measured both perpendicular and parallel to the axis of
 290 the core drilling forming two separate input parameters from the same petrophysical quantity. For the
 291 mercury injection and capillary pressure curves, one column of the data matrix contains those saturation
 292 values, which were measured on different core samples under the same pressure. The strength of
 293 correlation between the observed variables is found to be moderate. We introduce the following scalar
 294 called mean spread for the measure of average correlation:

$$295 \quad S = \left\{ \frac{1}{M(M-1)} \sum_{i=1}^M \sum_{j=1}^M [R(\mathbf{D})_{ij} - \delta_{ij}]^2 \right\}^{1/2}, \quad (7)$$

296 where $R(\mathbf{D})$ is the Pearson's correlation matrix of the observed variables, M is the number of
 297 measurement variables and δ is the Kronecker delta function. In this case study, S resulted around 0.6
 298 for the data matrix of core measurements. The linear regression connection between some of the studied
 299 petrophysical quantities is shown in Fig. 4, where the independent and dependent variables were denoted
 300 by x and y , respectively.

301

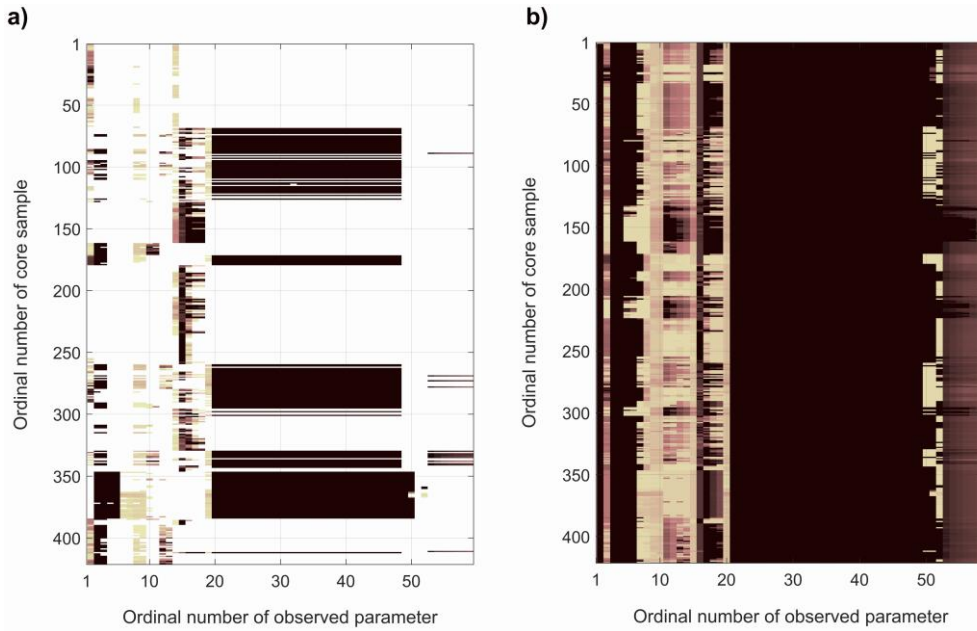


302

303 **Fig. 4.** Linear relationships between the petrophysical parameters in a Hungarian Miocene complex
 304 including the regression equations of porosity vs. bulk density (a), porosity measured perpendicular
 305 and parallel to the axis of the core drilling (b), Klinkenberg-corrected permeability vs. porosity (c) and
 306 water saturation measured by mercury injection and centrifugal methods (d)
 307

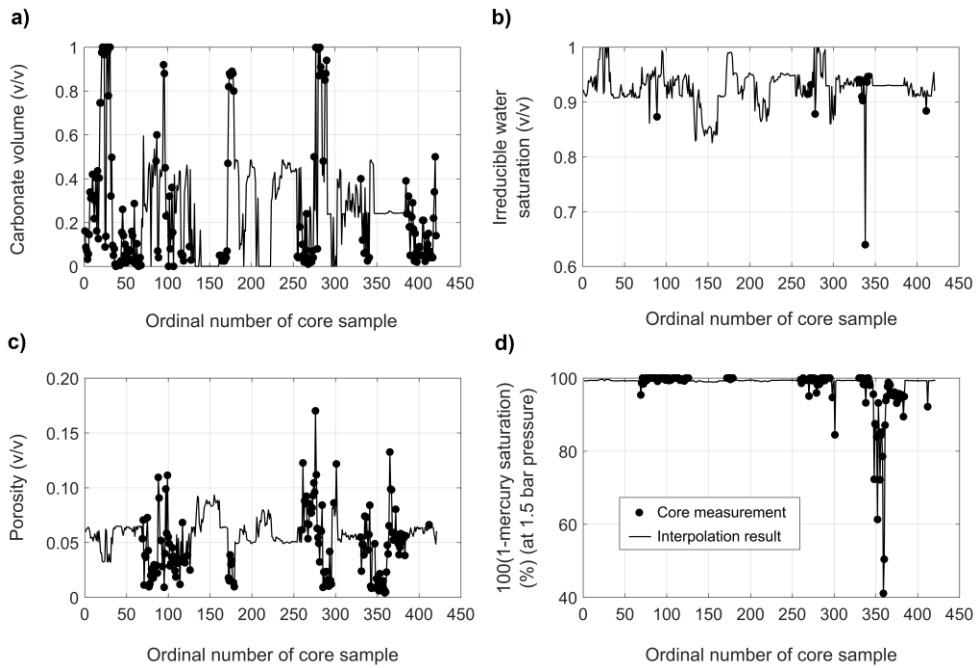
308 Core data referring to different wells and depth intervals are placed in the data matrix in random order.
 309 By using the linear regression functions determined between each pair of measured variables, the
 310 missing values in the input data matrix is replaced with synthetic data. Our correlation-based imputation
 311 method completely fills out the sparse matrix of the observed variables within their physical boundaries.
 312 The image of the original data matrix can be seen in Fig. 5a, where the missing data places are indicated
 313 by white color and the values of petrophysical parameters are scaled into the range of 0 and 1. The raw
 314 data matrix is incomplete at 71 %. One can notice that the data coming from different wells represent
 315 distinct groups of observed variables. It is shown in Fig. 5b that how the proposed statistical approach
 316 gives an estimate to the missing cells so that creating a complete data matrix of petrophysical quantities.
 317 At the end of this process, the large size data matrix is completed at 100 %. The result of interpolation
 318 is plotted for four essential petrophysical parameters in Fig. 6, which forms an appropriate input for the
 319 subsequent phase of the clustering procedure.

320



321

322 **Fig. 5.** Input data matrix includes missing values of petrophysical parameters in the investigated
 323 Hungarian Miocene complex (a) and fully filled data matrix extracted from the raw dataset using the
 324 correlation-based imputation method (b)
 325

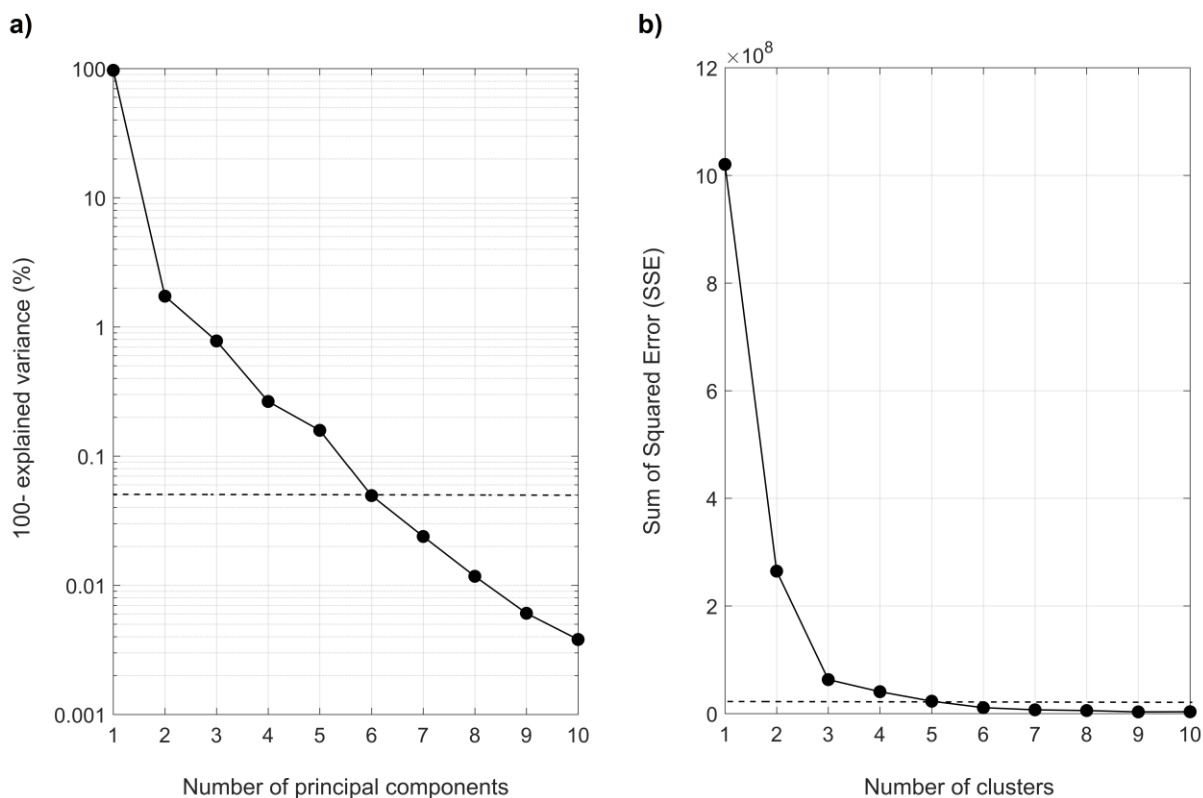


326

327 **Fig. 6.** Observed data measured on core samples (dots) and interpolation results for petrophysical
 328 parameters in a Hungarian Miocene complex (black line): carbonate content (a), irreducible water
 329 saturation (b), sample porosity (c), water saturation derived from mercury injection experiment (d)
 330

331 The resultant data matrix is decomposed by PCA according to Eq. (1), which concentrates the common
 332 observed information in a few variables. By calculating the eigenvalues of the PCs represented in Eq.
 333 (2), a decision can be made for the number of extracted variables. For the processed dataset, six PCs

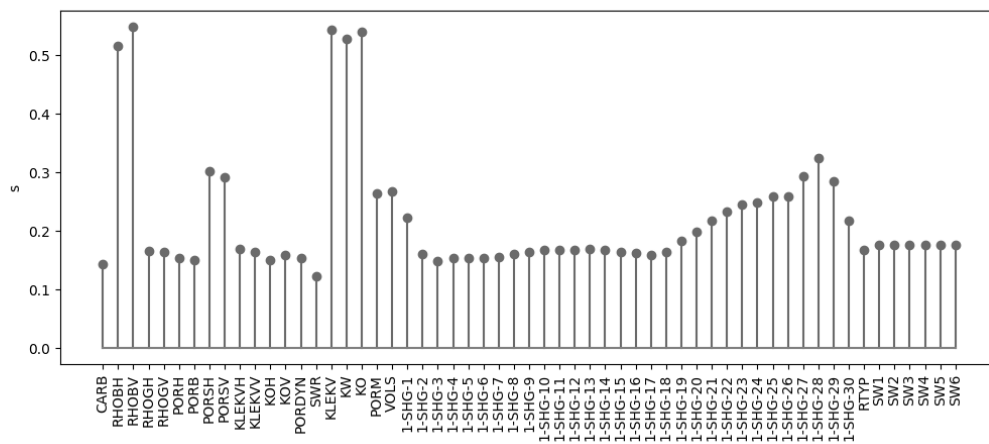
334 cover the 99.95 % of the total variance of the original data (Fig. 7a), so that the amount of information
 335 loss is 0.05 %. Then, the 421 objects in the six-dimensional space of PCs is clustered by a K-means
 336 method. The Elbow Method is a commonly used technique in the practice of cluster analysis to select
 337 the optimal number of clusters. Consequently, by using the SSE plot shown in Fig. 7b, five clusters were
 338 chosen by us for separating the main lithology types in the investigated Miocene complex.
 339



340
 341 **Fig. 7.** Variance portions explained by principal components derived from PCA (a), the sum of
 342 squared error calculated for different number of clusters given by K-means cluster analysis (b)
 343

344 The measured variables contribute to the solution to varying degrees. The PCA provides an opportunity
 345 for this analysis. The PC loadings given in Eq. (3) are plotted in Fig. 8. Based on the absolute values of
 346 the PCs' coefficients obtained in the input Miocene dataset, one can determine the relative weight of
 347 each petrophysical parameter. The figure shows the significant role of horizontal and vertical bulk
 348 density (RHOBH, RHOBV), horizontal and vertical porosity (PORSH, PORSV), Klinkenberg corrected
 349 permeability (KLEKV) and mercury injection data measured on higher pressures
 350 (1-SHG20,...,1-SHG28), while it indicates lower weights of the carbonate content (CARB), water

351 saturation (SW1,...,SW6) and mercury injection data at low pressures (1-SHG1,...,1-SHG17).
 352 Maximal impact on the PCs goes to the Klinkenberg corrected permeability and relative permeability to
 353 oil and water (KO, KW). In conclusion, PCA emphasizes primarily the pore structure characteristics
 354 and the hydraulic conductivity of the given rock types.

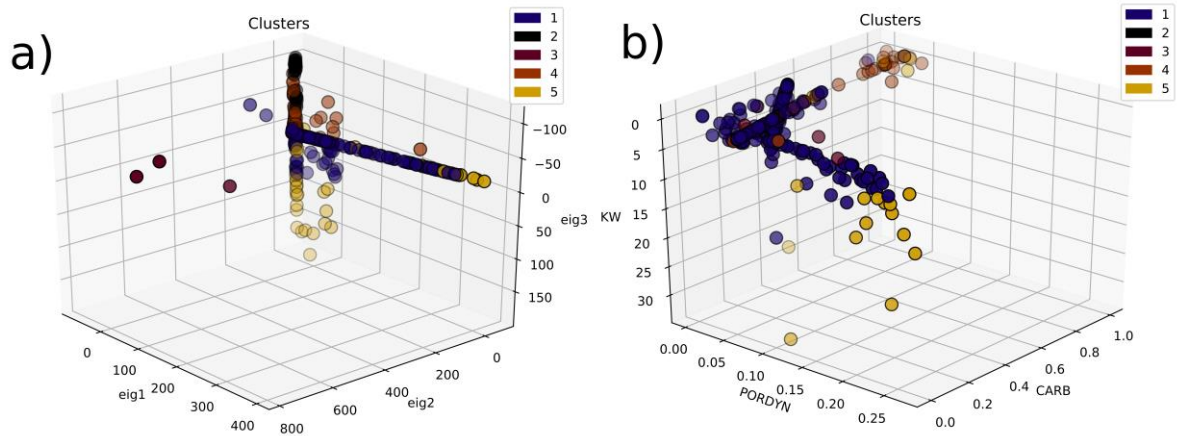


355

356 **Fig. 8.** Coefficients of principal components indicating the relative weights of petrophysical properties
 357 on the solution of PCA in the investigated Hungarian Miocene complex
 358

359 Clustering is made to separate the row vectors of observed petrophysical parameters into K=5 classes.
 360 To prevent the data processing from the harmful effect of outliers, the City block distance in Eq. (5) is
 361 used for the K-means cluster analysis. The grouped data objects in the space of the first three eigenvalues
 362 (eig1–eig3) is given in Fig. 9a, while the same objects are plotted in the coordinate system of three
 363 measured quantities (Fig. 9b) such as carbonate volume including calcite and dolomite contents
 364 (CARB), porosity obtained by dynamic displacement on full diameter samples (PORDYN) and relative
 365 permeability to water (KW). The figures demonstrate that the clusters are well separated (e.g., the high
 366 porosity and permeability reservoirs are indicated with dark yellow color) despite the input dataset being
 367 compositionally heterogeneous and unevenly distributed in the space.

368

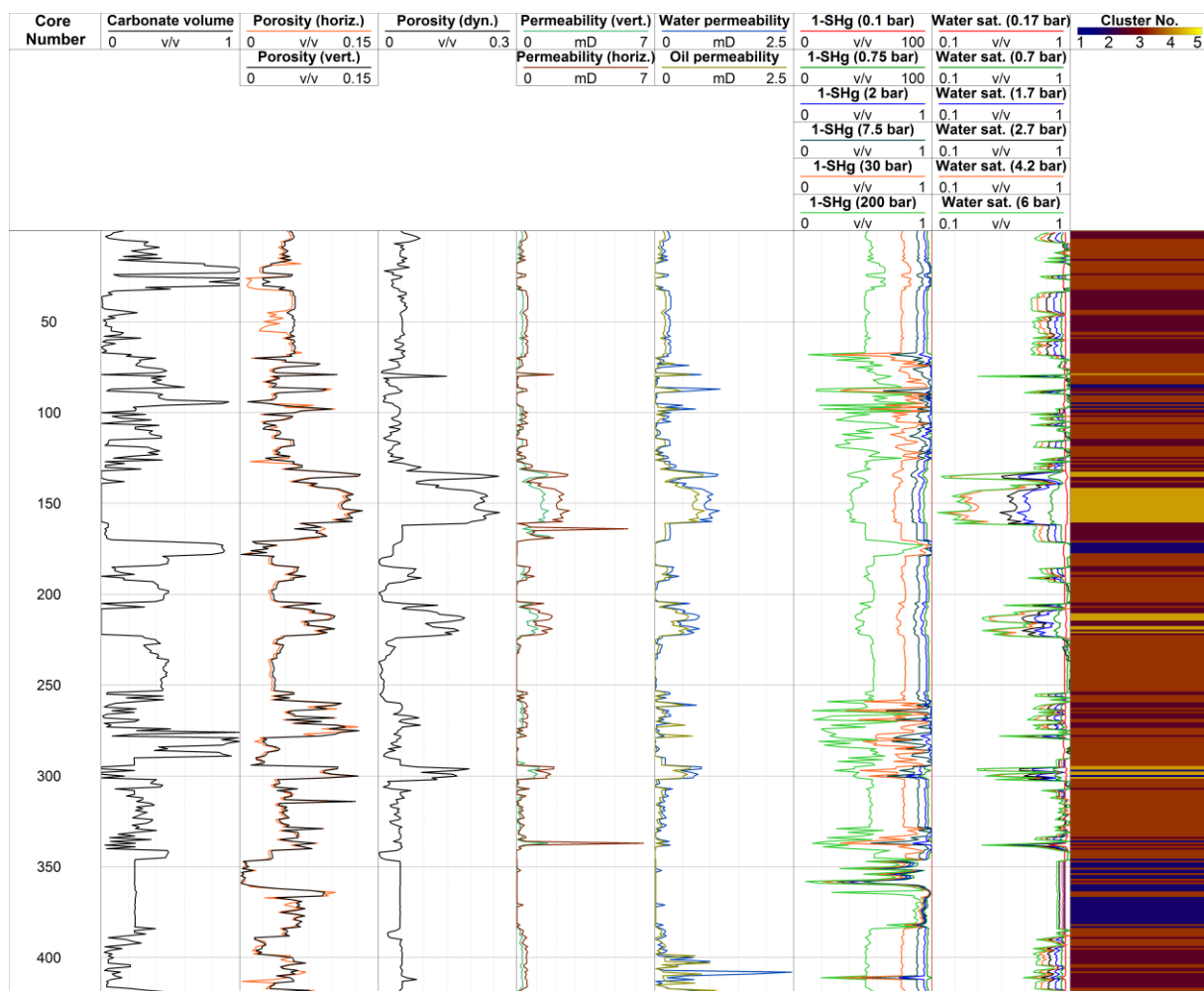


369

370 **Fig. 9.** Result of PCA of the filled data matrix observed in a Hungarian Miocene complex:
 371 transformed data in the coordinate system of eigenvectors (a). Result of cluster analysis: clusters vs.
 372 petrophysical variables as input parameters (b)
 373
 374

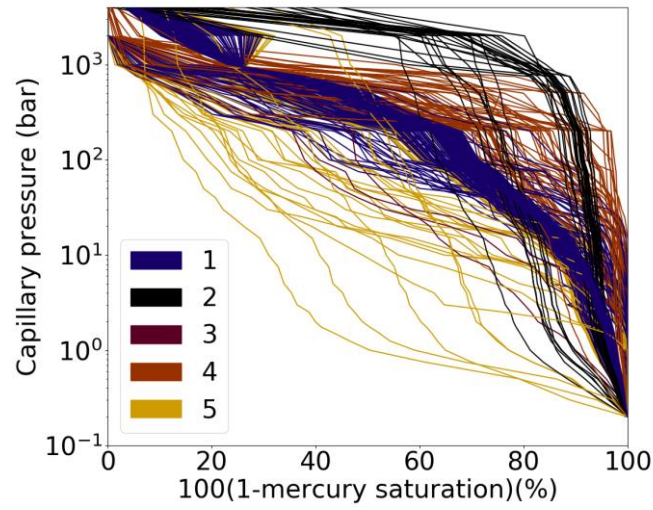
375 The result of cluster analysis can be seen in Fig. 10, where the section of cluster numbers in the last
 376 track shows to which group the rock samples are classified. It is noted that sudden changes in the
 377 resultant cluster numbers are also due to the fact that the core data referring to different depth coordinate
 378 (and wells) are quasi-randomly selected along the ordinate axis. The hydrocarbon-bearing formations
 379 can be clearly separated in the image of cluster numbers. Cluster 5 represented by yellow color is
 380 connected to highly porous and permeable zones with small amount of carbonate (e.g., for samples
 381 140–160). Petrophysical properties measured both with horizontal and vertical direction to the axis of
 382 sampling practically shows just a little amount of anisotropy. The capillary pressure curves indicates
 383 high amount of movable fluids in these intervals. The other four clusters indicate less permeable zones
 384 with higher amount of carbonate. Cluster 1 and 2 show impermeable formations with much amount of
 385 carbonate and high irreducible water saturation (e.g., around the near vicinity of sample 175), while
 386 cluster 3 and 4 form a transition between permeable and impermeable rocks (e.g., for samples 1–70).
 387 By comparing the result of cluster analysis to that of the lithology description, cluster 3 and 4 can be
 388 identified mostly as conglomerate and clay with some amount of dolomite, respectively. Cluster 5 is
 389 mainly composed of sandstone, gravel and fractured breccia with good reservoir storage capacity.
 390 Cluster 1 and 2 include calcareous and clayey marls. In the brown intervals, dolomites (e.g., around
 391 sample no. 200) and in the blue ones, thin layers of fractured basalt conglomerate are also found (e.g.,

392 in the range of 290–300). The clustering procedure allows the classification of capillary pressure curves
 393 by separating different groups (Fig. 11), which makes the interpretation more reliable. Although there
 394 is some overlap between the clusters of different pore geometries, the reservoirs of the highest movable
 395 wetting phase saturation can be clearly distinguished (see them by yellow color). The full curves of
 396 capillary pressure data gives more information than conventional porosity and permeability data. In
 397 addition to pore-size geometry distribution, it holds important information also on other textural
 398 characteristics of the clustered rock types, which can be further investigated in future studies.
 399



400

401 **Fig. 10.** Cluster analysis of laboratory data measured on core samples collected from a Hungarian
 402 Miocene complex: interpolated logs of petrophysical parameters as input for clustering (tracks 1-7),
 403 estimated distribution of clusters by using Manhattan distance metric as the result (track 8)
 404



405
406
407
408
409

Fig. 11. Results of clustering of mercury injection data measured on core samples collected from the studied Miocene formations

410 4. Discussion

411

412 To test the accuracy of estimation results of the proposed clustering method, a synthetic modeling
413 experiment has been accomplished. An exactly known inhomogeneous petrophysical model is assumed
414 to represent a hydrocarbon formation with varying amount of porosity (Φ), water saturation (S_w), sand
415 volume (V_{sd}), shale content (V_{sh}) and carbonate volume (V_c) along a borehole. Theoretical open-hole
416 wireline logging data are calculated and simultaneously processed by cluster analysis to examine how
417 accurately the exact model is reconstructed. The performance of cluster analysis can be tested using
418 arbitrarily chosen amount of noise added to the input data and the rate of incompleteness of the data
419 matrix formed from the noisy synthetic well logs. The following simplified probe response functions
420 summarized by Alberty and Hashmy (1984) are used to calculate the wireline logs (the effect of invasion
421 is neglected in this approach)

422

$$423 \quad \rho_b = \Phi [S_w \rho_w + (1 - S_w) \rho_h] + V_{sh} \rho_{sh} + V_{sd} \rho_{sd} + V_c \rho_c, \quad (8)$$

$$424 \quad GR = \rho_b^{-1} (V_{sh} GR_{sh} \rho_{sh} + V_{sd} GR_{sd} \rho_{sd} + V_c GR_c \rho_c), \quad (9)$$

425
$$P_e = \Phi [S_w P_{e,w} + (1 - S_w) P_{e,h}] + V_{sh} P_{e,sh} + V_{sd} P_{e,sd} + V_c P_{e,c}, \quad (10)$$

426
$$\Delta t = \Phi [S_w \Delta t_w + (1 - S_w) \Delta t_h] + V_{sh} \Delta t_{sh} + V_{sd} \Delta t_{sd} + V_c \Delta t_c, \quad (11)$$

427
$$\frac{1}{\sqrt{R_d}} = \left[\frac{V_{sh}^{\left(1 - \frac{V_{sh}}{2}\right)}}{\sqrt{R_{sh}}} + \frac{(\sqrt{\Phi})^m}{\sqrt{aR_w}} \right] (\sqrt{S_w})^n, \quad (12)$$

428 where the quasi-measured parameters are bulk density (ρ_b), natural gamma-ray intensity (GR),
 429 photoelectric absorption cross-section index (P_e), acoustic (P-wave) travel-time (Δt) and deep resistivity
 430 (R_d). The functional constants included in Eqs. (8)–(12), representing the physical properties of pore
 431 fluids, shale and mineral components as well as the textural properties of rocks, are specified in Table 1
 432 (we assume gas phase under hydrocarbon).

433
 434 **Table 1**
 435 Zone parameters used for calculating synthetic wireline logs to test the accuracy of the proposed cluster
 436 analysis method.
 437

Well log	Zone parameter	Symbol	Selected value	Unit
Natural gamma intensity (GR)	sand	GR _{sd}	10	API
	shale	GR _{sh}	140	
	carbonate	GR _c	5	
Bulk density (ρ_b)	sand	ρ_{sd}	2.65	g/cm ³
	shale	ρ_{sh}	2.55	
	carbonate	ρ_c	2.79	
	pore-water	ρ_w	1.09	
	hydrocarbon	ρ_h	0.016	
Sonic interval-time (Δt)	sand	Δt_{sd}	56	$\mu\text{s}/\text{ft}$
	shale	Δt_{sh}	108	
	carbonate	Δt_c	46	
	pore-water	Δt_w	200	
	hydrocarbon	Δt_h	305	
Photoelectric index (P_e)	sand	$P_{e,sd}$	1.81	barn/e
	shale	$P_{e,sh}$	3.50	
	carbonate	$P_{e,c}$	4.11	
	pore-water	$P_{e,w}$	0.81	
	hydrocarbon	$P_{e,h}$	0.09	
Deep resistivity (R_d)	shale	R_{sh}	1	Ωm
	pore-water	R_w	0.06	
	cementation exponent	m	2.0	–
	saturation exponent	n	2.0	
	tortuosity coefficient	a	1.0	

438 In this uncertainty analysis, 5 % Gaussian distributed noise is added to the synthetic data calculated in
 439 the forward modeling procedure using Eq. (8)-(12). The mean spread (defined in Eq. (7)) calculated for
 440 the noisy well logging parameters is 0.44, which represents weak correlation relationship between the
 441 five input variables on the average. The full correlation matrix is given in Table 2.

442

443 **Table 2**

444 Pearson's correlation matrix of synthetic well logs including 5 % Gaussian noise.

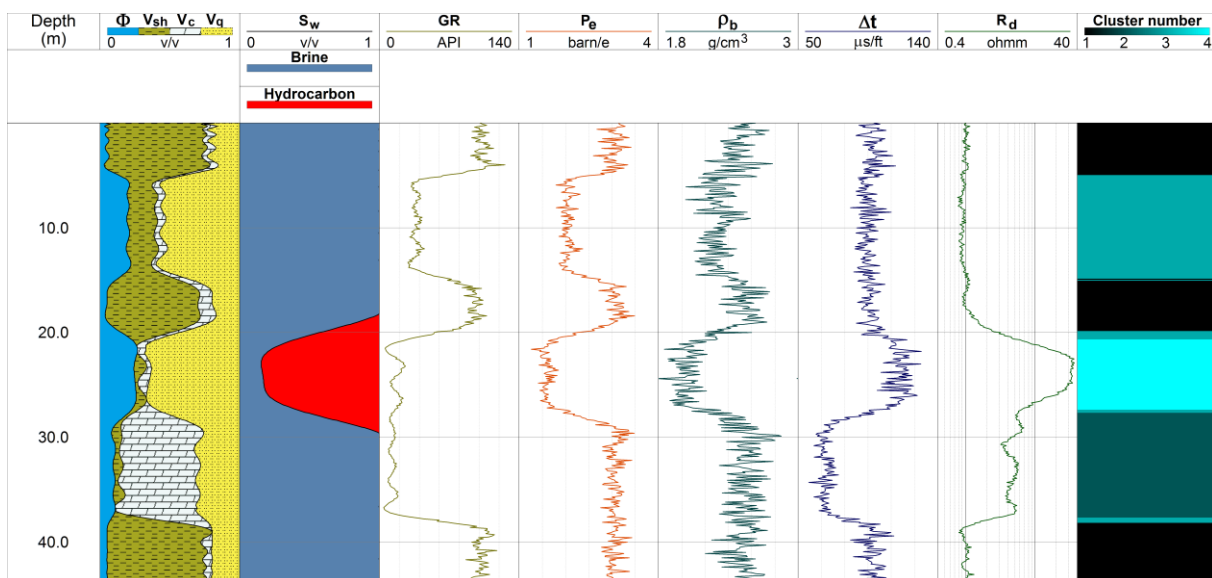
445

	GR	P_e	ρ_b	Δt	R_d
GR	1	0.54	0.44	0.23	-0.50
P_e	0.54	1	0.61	-0.38	-0.49
ρ_b	0.44	0.61	1	-0.34	-0.50
Δt	0.23	-0.38	-0.34	1	0.25
R_d	-0.50	-0.49	-0.50	0.25	1

446

447 The exactly known petrophysical model can be seen in tracks 1–2 in Fig. 12, where 6 layers can be
 448 distinguished: 5 m thick shale, 10 m thick water-bearing sandstone, 5 m thick shale, 8 m thick
 449 hydrocarbon-bearing sandstone, 10 m thick limestone and 6 m thick shale. By assuming 4 main
 450 lithological units, the entire well logging dataset (plotted in tracks 3-7) is processed by K-means cluster
 451 analysis (in this test PCA is not necessary to be applied). As a result, the depth variation of cluster
 452 numbers is shown in the last track.

453

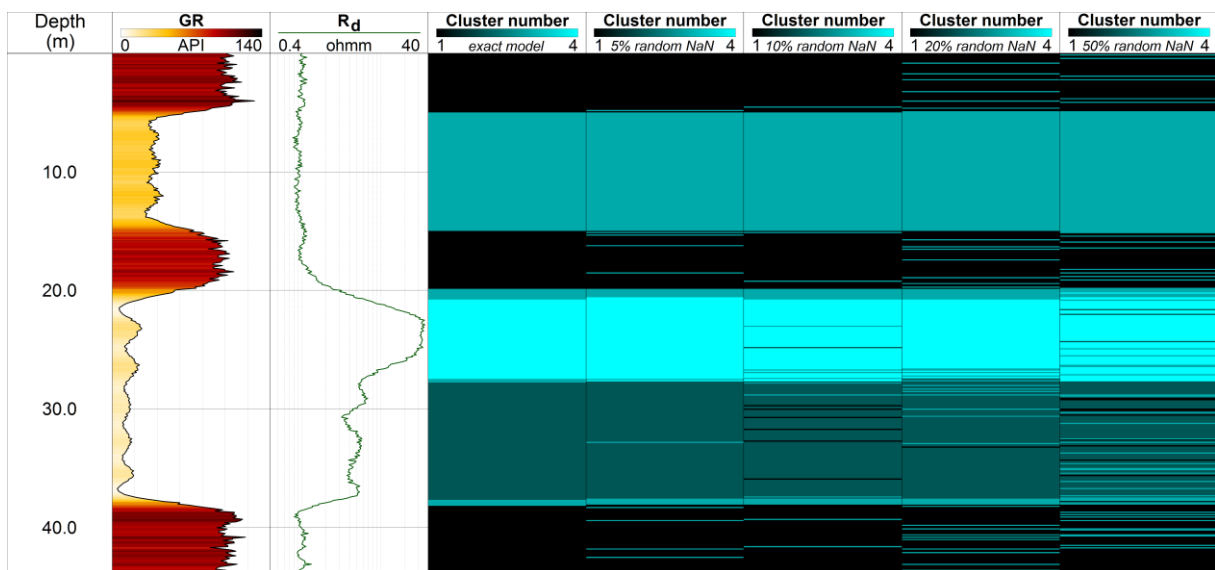


454

455 **Fig. 12.** Cluster analysis of synthetic wireline logs contaminated with 5 % Gaussian distributed
 456 random noise
 457

458 The reliability of the correlation-based imputation method assisted cluster analysis is tested using
 459 different datasets derived from the noisy synthetic well logs, where the proportion of missing data is
 460 gradually increased. The vertical distribution of the four clusters determined from the noiseless data is
 461 illustrated on the third track of Fig. 13, which is considered to be the exact model. For this analysis, a
 462 number of data randomly chosen from all logs are substituted with NaN values to form new input
 463 datasets and correlation-based imputation is performed before clustering them. To the right from the
 464 third track, the result of cluster analysis is shown for the incomplete well logging datasets, where the
 465 number of NaN values relative to the total number of data is 5, 10, 20 and 50 %, respectively. The figure
 466 shows that even in large incompleteness of the input data matrix the main lithological units can be
 467 reliably separated. These results are properly good for a relatively short depth interval (and few number
 468 of data) and small number of weakly correlated measured variables. The uncertainty of estimation
 469 depending on the accuracy of data and interpolation errors can be reduced by selecting input variables
 470 with stronger correlations and removing the outliers from the dataset. The latter can be effectively done
 471 by using automated robust clustering methods.

472



473

474 **Fig. 13.** Effect of the amount of data matrix incompleteness on the result of cluster analysis
 475
 476

477 The aforementioned synthetic modeling experiment focuses mainly on the effect of the accuracy of
478 original measurements and the interpolation error when replacing the unknown values with estimated
479 ones. Besides these types of noises added to the dataset, the data also include other sources of
480 uncertainty. For a more detailed uncertainty analysis, we have to consider basically three types of errors.
481 At first, one has to quantify the error of the original measurements. The estimation of measurement
482 errors makes it necessary to know the features of the measurement process. Secondly, our model
483 assumes linear connection between all the columns in the data matrix. For a better approximation, one
484 has to consider the underlying physical relations among the columns, which is unknown in most of the
485 cases. At last, the accumulated estimation errors should also be taken into consideration. It is the result
486 of the incremental nature of the imputation algorithm giving an estimate in the actual iteration from the
487 estimated values of earlier iterations. In the ideal case, only the considered value is missing, which
488 means that the estimation is based on the available measurements. In the worst case, all values of the
489 selected columns are estimated. This kind of accumulated estimation error depends on the rate of the
490 counts of measured and predicted values having used for the given estimation.

491

492 **5. Conclusions**

493

494 In the paper, an improved clustering algorithm is proposed for the interpretation of core datasets
495 originating from heterogeneous sources. A correlation-based multi-linear regression method is first used
496 for the filling of typically incomplete data matrices, which gives a suitable input for the K-means cluster
497 analysis performed in a subsequent data processing step. After the upload phase, all variables are taken
498 into account during the clustering, but it is also possible to select the parameters individually based on
499 the results of principal component analysis. Those parameters with relative small PC coefficients can be
500 neglected. Since there is no limit for the approach on the number of processed variables forming the
501 input dataset, the clustering method is suitable for the future expansion of observed parameters. Based
502 on this, not only new petrophysical parameters, but also mineralogical (e.g., from X-ray diffraction) and
503 other compositional (e.g., geochemistry) data may be involved. An added advantage of the statistical
504 method is the grouping possibility of the capillary pressure curves, which helps to integrate more

505 information on the textural and fluid characteristics of the investigated formations. (The analysis of these
506 curves has wide literature, which is beyond the scope of this paper.) The explanation of the behavior of
507 these curves are typically given on an empirical basis in the literature. However, an analytic description
508 may reveal new petrophysical characteristics of the studied formations. In addition to the measured
509 parameters of such analytic models, for a more reliable interpretation of the pore structure (textural
510 properties), pore-fluid types (bound or free water and hydrocarbons) and properties (e.g., viscosity), T2
511 relaxation time distribution curves of nuclear magnetic resonance measurements can also be added as
512 input in future applications.

513

514 Cluster analysis is commonly used also for the processing of wireline logging data. Well log derived
515 petrophysical parameters such as porosity, shale volume, matrix volumes, water and hydrocarbon
516 saturation, permeability etc. as high resolution in situ information can significantly increase the size of
517 the statistical sample and may further improve the performance of cluster analysis. The big amount of
518 data is processed by linear regression and classical K-means clustering algorithm for a quick
519 interpretation. For obtaining a more reliable solution, a robustified cluster analysis approach including
520 a nonlinear interpolation method as imputation phase can be used in the future. As a statistically highly
521 efficient method, the most frequent value-based (automated) weighting procedure introduced by Steiner
522 (1991) gives a robust solution independent from the nature and statistical distribution of the input dataset
523 (e.g., in Zhang, 2017). As a first result, its application to clustering of well logs can be found in Braun
524 et al. (2016), which will be followed by improved clustering algorithms in the near future as intensive
525 research is currently made at the Department of Geophysics, University of Miskolc. The resultant
526 clusters can also be easily correlated between the wells to represent the spatial distribution of clusters
527 for two- and three-dimensional cases.

528

529 **Acknowledgments**

530

531 The authors would like to thank the experts and staff at the MOL group for their valuable cooperation
532 and for the measurement and laboratory data they provided. This research was supported by the

533 European Union and the Hungarian State, co-financed by the European Regional Development Fund in
534 the framework of the GINOP-2.3.4-15-2016-00004 project, aimed to promote the cooperation between
535 the higher education and the industry.

536

537 **References**

538

539 Alberty, M., and K. Hashmy, 1984. Application of ULTRA to log analysis: Presented at the SPWLA
540 Symposium Transactions, 1–17.

541 Braun B. A., Abordán A., Szabó N. P., 2016. Lithology determination in a coal exploration drillhole
542 using Steiner weighted cluster analysis. *Geosciences and Engineering* 5 (8), 51–64.

543 Candès E. J., Recht B., 2009. Exact matrix completion via convex optimization. *Foundations of*
544 *Computational mathematics* 9(6), 717–772.

545 Cranganu C., Luchian H., Breaban M. E., 2015. Artificial intelligent approaches in petroleum
546 geosciences. Springer.

547 Dolton G. L., 2006. Pannonian Basin Province, Central Europe (Province 4808) – Petroleum geology,
548 total petroleum systems, and petroleum resource assessment, *USGS Bull.* 2204–B, 1–47.

549 Gemulla R., Nijkamp E., Haas P. J., Sismanis Y., 2011. Large-scale matrix factorization with distributed
550 stochastic gradient descent. In Proceedings of the 17th ACM SIGKDD international conference on
551 Knowledge discovery and data mining, ACM, 69–77.

552 Hartigan J. A., Wong M. A., 1979. Algorithm AS 136: A k-means clustering algorithm. *Journal of the*
553 *Royal Statistical Society. Series C (Applied Statistics)*, 28.1, 100–108.

554 Hempkins W. B., 1978. Multivariate statistical analysis in formation evaluation. SPE California
555 Regional Meeting, San Francisco, 7144–MS.

556 Jarzyna J. A., Bała M., Krakowska P. I., Puskarczyk E., Strzępowicz A., Wawrzyniak-Guz K., Więclaw
557 D., Ziętek J., 2017. Shale Gas in Poland, *Advances in Natural Gas Emerging Technologies*. IntechOpen,
558 DOI: 10.5772/67301.

559 Jolliffe I.T., 2002. *Principal component analysis*, 2nd edition, New York, Springer-Verlag.

560 Jung H., Jo H., Kim S., Lee K., Choe J., 2018. Geological model sampling using PCA-assisted support
561 vector machine for reliable channel reservoir characterization. *Journal of Petroleum Science and*
562 *Engineering* 167, 396–405.

563 Lawley D. N., Maxwell A. E., 1962. Factor analysis as a statistical method. *The Statistician* 12, 209–
564 229.

565 Little R., Rubin B., 1986. *Statistical Analysis with Missing Data*, John Wiley & Sons, Inc., New York.

566 Ma Z., Holditch S., 2015. *Unconventional Oil and Gas Resources Handbook: Evaluation and*
567 *Development*. Gulf Professional Publishing.

568 Masoudi P., Aifa T., Memarian H., Tokhmechi B., 2018. Uncertainty assessment of porosity and
569 permeability by clustering algorithm and fuzzy arithmetic. *Journal of Petroleum Science and*
570 *Engineering* 161, 275–290.

571 Mazumder R., Hastie T., Tibshirani R., 2010. Spectral regularization algorithms for learning large
572 incomplete matrices. *Journal of Machine Learning Research* 11, 2287–2322.

573 Skalinski M., Gottlib-Zeh S., Moss B., 2006. Defining and predicting rock types in carbonates –
574 Preliminary results from an integrated approach using core and log data from the Tengiz Field.
575 *Petrophysics* 47 (1), 37–52.

576 Steiner F., 1991. *The most frequent value: Introduction to a modern conception of statistics*. Akadémiai
577 Kiadó, Budapest.

578 Szabó N. P., Dobróka M., Kavanda R., 2013. Cluster analysis assisted float-encoded genetic algorithm
579 for a more automated characterization of hydrocarbon reservoirs. *Intelligent Control and Automation* 4
580 (4), 362–370.

581 Szabó N. P., Dobróka M., 2017. Robust estimation of reservoir shaliness by iteratively reweighted factor
582 analysis. *Geophysics* 82 (2), D69–D83.

583 Szabó N. P., Dobróka M., 2018. Exploratory factor analysis of wireline logs using a Float-Encoded
584 Genetic Algorithm. *Mathematical Geosciences* 50 (3), 317–335.

585 Zhang J., 2017. Most frequent value statistics and distribution of ^7Li abundance observations. *Monthly*
586 *Notices of the Royal Astronomical Society* 468 (4), 5014–5019.

587

588 **Appendix**

589

590 A simplified pseudo-code of the imputation algorithm applied in the proposed workflow, where m
591 denotes the total number of columns of the input data matrix (i.e., measured petrophysical parameters,
592 here variables x and y) and n is the total number of rows of the same matrix (i.e., core samples).

593

```
594 for i := 1 to m - 1
595     for j := i + 1 to m
596         x := i-th column of the data matrix
597         y := j-th column of the data matrix
598         points = {(xk, yk) : xk != NaN, yk != NaN, 1 <= k <= n}
599         fit linear model to points by linear regression
600     for k := 1 to n
601         if xk != NaN and yk == NaN
602             yk := fitted value by the linear model from xk
603         end
604     if xk == NaN and yk != NaN
605         xk := fitted value by the linear model from yk
606     end
```

607 end
608 end
609 end