# Machine learning approach for predicting crude oil stability based on NMR spectroscopy

Dubravka Raljević[a], Jelena Parlov Vuković[a,*], Vilko Smrečki[b,*], Ljiljana Marinić Pajc[a],

Predrag Novak[c,*], Tomica Hrenar[c,*], Tomislav Jednačak[c,*], Lucija Konjević[a], Bruno Pinević[c],

Tonka Gašparac[c]

[a] *INA-Industrija nafte d.d., Refining & Marketing, Central Testing Laboratory, Lovinčićeva 4, HR-10002 Zagreb, Croatia*

[b] *Ruđer Bošković Institute, NMR Centre, Bijenička 54, HR-10000 Zagreb, Croatia*

[c] *University of Zagreb, Faculty of Science, Department of Chemistry, Horvatovac 102a, HR-10000, Zagreb, Croatia*


\* Corresponding authors.

*E-mail addresses:* jelena.parlov-vukovic@ina.hr (J.P.V.); smrecki@irb.hr (V.S.); pnovak@chem.pmf.hr (P.N.); hrenar@chem.pmf.hr (T.H.); tjednacak@chem.pmf.hr (T.J.)

**Keywords:** machine learning, crude oil, NMR spectroscopy, stability

**Abstract**

Crude oils are extremely complex organic mixtures, composed of various constituents ranging in size, shape and polarity. Obtaining a detailed insight into the petroleum composition is of highest priority for quality evaluation of crude oils and crude oil product performances. The stability of crude oils and their components represents one of the major challenges in petroleum industry, since there is no existing single method to determine the stability of all fractions. In this study, statistical multi-way analysis (MWA) and machine learning (ML) methods were coupled with diffusion-ordered NMR spectroscopy (DOSY) and compared to different crude oil stability affecting parameters in order to explore possibilities to predict crude oil stability. The potential of this approach was explored to identify and classify the crude oils of different origin according to their composition, stability, density and diffusion properties. With the application of MWA using the TUCKER3 decomposition model for a set of DOSY NMR spectra, the principal components were determined for the model (5,5,5), which described 99.89% of the total variance. The reduced space of the first 3 principal components was used for the sample classification. Similar samples were identified, and reduced space was further utilized for the regression of measured stabilities. Extensive ML multivariate linear regression was carried out for modeling crude oil stability in relation to DOSY NMR spectra and other measured properties, such as aromaticity, API gravity, percentage of aliphatic chains, asphaltene content and relative diffusivities. In both MWA and ML cases the best predictive models were determined. For such complex mixtures as crude oils are, exceptionally good correlations were obtained, proving that this new and robust model can accurately predict crude oil stability and other important parameters relevant for petroleum industry thus showing a great potential for practical applications.

## 1. Introduction

Crude oil is a highly complex organic mixture composed of various aliphatic and aromatic hydrocarbons ranging in size, shape and polarity. Obtaining a detailed insight into chemical composition is of highest priority for quality evaluation of crude oils and crude oil product performances [1,2]. According to their polarity, the components of crude oil are often divided into four main groups: asphaltenes, saturates, aromatics and resins. Asphaltenes are the heaviest and the most polar crude oil components, composed of aromatic and saturated rings, aliphatic moieties, some heteroatoms, such as nitrogen, oxygen and sulfur, and traces of transition metals [3–5]. During petroleum processing, asphaltenes may form aggregates and precipitates, leading to serious problems in production, transportation and storage. The stability of asphaltenes in crude oils and petroleum products is one of the major challenges in petroleum industry, since there is no existing single method to determine stability of all oil fractions [6–9].

Nuclear magnetic resonance (NMR) spectroscopy has emerged as a valuable tool for studying crude oils and their derivatives [10–18]. However, proton and carbon NMR spectra of petroleum samples are characterized by severely overlapping signals, which are difficult to straightforwardly assign and analyze. Further insight into the nature and structure of crude oils can be obtained by diffusion ordered NMR spectroscopy (DOSY) [14–18]. This approach can be applied to measure translational diffusion properties of individual components in complex mixtures without their physical separation. DOSY NMR spectra are pseudo-two-dimensional, where one dimension is represented by chemical shifts and the other by translational diffusion coefficients, which depend on the shape and size of a molecule or an aggregate in the sample. Crude oils originating from different geographical regions contain various types of compounds that can be separated and identified according to their diffusion coefficients. Nevertheless, even with the state-of-the-art NMR techniques one is still not able to perform a complete

71  differentiation among crude oil samples based on spectral inspection only (Figs. 2 and 3).

72  Hence, further evaluation and spectral processing by statistical methods are required to explore

73  the correlation between the origin and physical properties of crude oils. Recently, it has been

74  shown that petroleum samples of different origin can be identified, clustered and well-separated

75  by employing a combination of DOSY NMR spectroscopy and multi-way analysis [17].

76  Moreover, an advanced statistical model based on trilinear decomposition algorithm has been

77  developed, validated and applied to evaluate DOSY NMR spectra. In a similar study, proton

78  NMR spectra have been processed by principal component analysis to reveal characteristic

79  spectral areas responsible for sample differentiation and classification [18].

80      In this study, multi-way analysis and machine learning methods are combined to predict

81  the crude oil stability based on NMR spectroscopy. For this purpose, capabilities of DOSY

82  NMR coupled with both multi-way and machine learning multivariate linear regression

83  analyses have been explored to identify and classify crude oils of different origin according to

84  their content, stability, density and diffusion properties.

85  **2.  Experimental**

86  *Samples*

87      All crude oil samples were obtained from geographical regions designated in Tables 1,

88  2 and S1.

89  *2.1.  Asphaltene content analysis*

90      Asphaltenes were extracted from the crude oil samples by employing the standard

91  ASTM D 6560-17 method to determine the content of insoluble asphaltenes in heptane [19].

92  The crude oil samples were refluxed in heptane and mixed with the precipitate. Subsequently,

93  asphaltenes, waxy substances and inorganic material were collected on a filter paper. In the

94  next step, the waxy substances were removed by washing with hot heptane in an extractor,

95  while the asphaltenes were separated from the inorganic material by dissolving in hot toluene.

96  The extraction solvent was evaporated.

97  *2.2.  Stability testing*

98      Stability testing was based on the ASTM D 7157-18 standard method [20]. The sample

99  solutions were prepared in toluene at three different concentrations and analyzed by ROFA

100 France automated stability analyzer equipped with an optical probe for detecting the asphaltene

101 flocculation. Stability parameters $S_{total}$ (overall stability of the sample), $S_{asph}$ (peptizability or

102 ability of asphaltenes to remain in a dispersed state) and $S_{resin}$ (aromaticity of the resins and their

103 capability to maintain asphaltenes in solution) were calculated as well as intrinsic stability from

104 volumes of toluene and *n*-heptane and mass of the samples.

105 *2.3.  NMR measurements*

106     NMR experiments were performed at 298 K and chemical shifts were reported relative

107 to tetramethylsilane (TMS) internal standard. The samples (100 μL) were dissolved in 500 μL

108 of a deuterated solvent. One-dimensional $^1$H NMR spectra were recorded on a Bruker Avance

109 Neo 300 NMR spectrometer in chloroform-d (99.8%, Aldrich) using a C/H dual 5 mm probe

110 with 32 scans, 10 s recycle delay, 7.6 μs $\pi/2$ pulse length and 16 K time domain. $^1$H NMR

111 DOSY spectra were acquired in toluene-$d_8$ (99.5%, CIL) on a Bruker Avance 600 NMR

112 spectrometer using a 5 mm TBI probe equipped with z-gradients. Triplicate measurements were

113 carried out using a dstebpgp3s pulse sequence with convection compensation, 16 scans,

114 6.0 kHz spectral width, 600 μs spoil gradients, 16 K time domain, 150 μs gradient recovery and

115 5 ms eddy current delays. The gradient strength was varied from 2 % to 95 % in 16 steps, while

116 both the small (2.0 ms) and the big (70.0 ms) delta were kept constant.

117 *2.4.  Multi-way analysis*

118     *Multi-way* analysis (MWA) presents decomposition of multidimensional datasets

119 represented as multidimensional numerical arrays (or a higher order data tensor). It could be
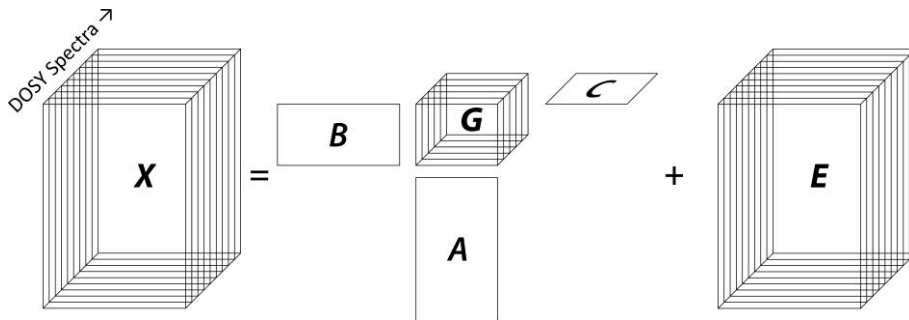
120    considered as an extension of principal component analysis [21]. Data tensor is composed from

121    sequences of numbers dependent on different physical dimensions or *ways*. In this case, the 3$^{rd}$

122    order tensor consists of two-dimensional DOSY NMR spectra for different crude oil samples.

123        Each DOSY NMR spectrum was extracted with 128×2192 records providing the total

124    dimensions of the 3$^{rd}$ order tensor: 18×128×2192. The data in this 3$^{rd}$ order tensor depend on

125    three independent variables: chemical shift, magnetic gradient pulse amplitude and sample

126    diversity [17]. To extract the quantitative classification information, MWA was used as a tool

127    that allows detection of variabilities among all investigated samples based on their

128    2-dimensional DOSY NMR datasets. After tensor decomposition, each DOSY NMR spectrum

129    was finally represented as one point in reduced space.

130        MWA on the set of DOSY NMR spectra placed in the 3$^{rd}$ order tensor was carried out

131    using the 3-way decomposition model TUCKER3 [22]:

132    $$\boldsymbol{X} = \boldsymbol{AG}(\boldsymbol{C} \otimes \boldsymbol{B})^{\tau} + \boldsymbol{E} \tag{1}$$

133    where $\boldsymbol{A}$, $\boldsymbol{B}$, and $\boldsymbol{C}$ are the 1$^{st}$-way, 2$^{nd}$-way, and 3$^{rd}$-way loadings matrices, respectively (symbol

134    $\otimes$ represents Kronecker matrix product) (Fig. 1).



**Fig. 1.** Graphical representation of the TUCKER3 model.

137        The $\boldsymbol{G}$ matrix is the *core-array* and is associated with the amount of variation explained

138    by loadings in the different modes. MWA was performed by using the code ***moonee*** [23–26]

139    developed *in-house*.

6

140   *2.5.   Machine learning multivariate linear regression*

141        Measured stability data was used as dependent variable in two cases. Firstly, this data

142   was regressed on the first three principal components of the reduced space of DOSY NMR

143   spectra. In the second case, stability data was regressed to 6 other measured properties.

144   Extensive machine learning (ML) procedure was applied for generation of all possible

145   multivariate linear regression (MLR) models with any possible linear combination of original

146   variables as well as their higher-order polynomial terms (up to the $4^{th}$ order in the first case and

147   up to the $2^{nd}$ order in the second case). Total numbers of generated different models for each

148   dependent variable were 1 717 869 184 and 134 217 728 for the first and the second case,

149   respectively.

150        MLR was performed using the following expression for matrices of coefficients B

151   calculated by singular value decomposition:

$$B = (X^\tau X)^{-1} X^\tau Y \tag{2}$$

152

153   where $X$ and $Y$ are the matrices of independent and dependent variables, respectively. For each

154   model, *leave-one-out cross-validation* (LOO-CV) was performed, and various statistical

155   parameters were computed. Methodical validation of models by LOO-CV provided an optimal

156   representation selected on the basis of adjusted and predicted $R^2$ values as well as LOO-CV

157   mean squared error.

158   **3.    Results and Discussion**

159   *3.1.   Evaluation of crude oil properties by standard testing methods*

160        Chosen crude oil properties, such as the asphaltene content ($w_{asph}$), stability parameters

161   ($S_{asph}$, $S_{resin}$, $S_{total}$) and API gravity values ($\rho_{API}$) of analyzed crude oil samples are shown in

162   Table 1. API values indicated that all analyzed crude oils belong to lighter and medium crude

163   oil categories.

164       The overall stability of crude oils decreases with the increase in asphaltene content.

165 However, a comparison of some test samples show that the crude oil stability does not

166 exclusively depend on the content of asphaltenes, especially in light crude oils. Parameters,

167 such as the composition and structure of resins, aromatics and other components affect the total

168 stability. Furthermore, the stability of crude oils depends on various processes that involve

169 blending, dilution, temperature and pressure changes.

170 **Table 1.** Comparison of crude oil properties determined by standard methods.

| Sample No. | Designation | $w_{asph}$ / % [a] | $S_{asph}$ [b] | $S_{resin}$ [c] | $S_{total}$ [d] | $\rho_{API}$ [e] |
|---|---|---|---|---|---|---|
| 1 | North Africa 1 | 2.50 | 0.75 | 0.58 | 2.33 | 29.93 |
| 2 | Southwest Asia 1 | 0.15 | 0.41 | 1.88 | 3.18 | 38.15 |
| 3 | Southwest Asia 2 | 2.25 | 0.6 | 1.15 | 2.92 | 31.21 |
| 4 | Southwest Asia 3 | 2.74 | 0.76 | 0.57 | 2.39 | 30.13 |
| 5 | Central Europe 1 | 1.19 | 0.77 | 0.69 | 3.02 | 30.11 |
| 6 | Eastern Europe 1 | 0.80 | 0.77 | 1.07 | 4.61 | 29.58 |
| 7 | Eastern Europe 2 | 0.26 | 0.78 | 1.01 | 4.54 | 36.39 |
| 8 | Central Europe 2 | 1.01 | 0.77 | 0.73 | 3.22 | 30.24 |
| 9 | West Africa | 0.01 | 0.44 | 1.63 | 2.94 | 32.42 |
| 10 | Southwest Asia 4 | 3.12 | 0.74 | 0.59 | 2.25 | 29.94 |
| 11 | Southwest Asia 5 | 1.58 | 0.77 | 0.48 | 2.12 | 35.24 |
| 12 | Eastern Europe 3 | 0.64 | 0.65 | 0.95 | 2.72 | 35.05 |
| 13 | NorthEast Asia | 1.45 | 0.77 | 0.68 | 2.97 | 29.58 |
| 14 | North Africa 2 | 0.37 | 0.65 | 0.81 | 2.33 | 37.15 |
| 15 | North Asia | 0.71 | 0.75 | 1.15 | 4.68 | 33.51 |
| 16 | Central Europe 3 | 0.14 | 0.64 | 0.97 | 2.71 | 37.51 |
| 17 | Central America | 0.62 | 0.70 | 1.43 | 4.69 | 41.70 |

171 [a]   asphaltene content according to gravimetric analysis
172 [b]   peptizability or ability of asphaltenes to remain in a dispersed state
173 [c]   aromaticity of the resins and their capability to maintain asphaltenes in solution
174 [d]   total stability or overall stability of the sample
175 [e]   gravity according to American Petroleum Institute (API)

176
177 *3.2.  NMR spectroscopy*

178       Typical proton NMR spectra of crude oil samples are displayed in Fig 2. Severe peak

179 overlapping makes these spectra difficult to analyze and only information on different classes
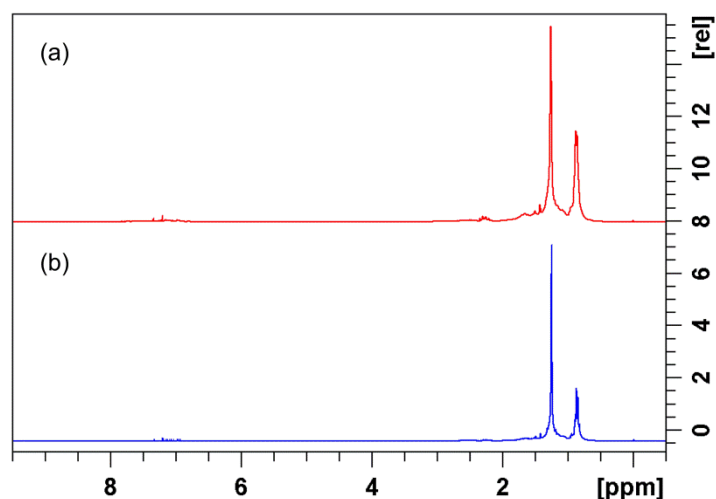
8

180 of hydrocarbons can be obtained. One of the features that can be determined from $^1$H NMR

181 spectra is the aromaticity ($H_{ar}$), usually expressed as the content (in percentage) of aromatic

182 hydrogen atoms. It can be calculated as the ratio between the sum of all aromatic hydrogen

183 integrals ($I_{H_{ar}}$) and the total amount of hydrogen atoms (consisting of the sum of all aliphatic

184 and aromatic hydrogen integrals, $I_{H_{aliph}}$ and $I_{H_{ar}}$), using the previously described procedure

185 [18,27,28]:

$$H_{ar}[\%] = \frac{\sum I_{H_{ar}}}{\sum I_{H_{aliph}} + \sum I_{H_{ar}}} \qquad (3)$$

187       In order to assure more accurate $H_{ar}$ calculation by avoiding overlapping signals of crude

188 oil aromatic hydrogens with toluene aromatic hydrogens, corresponding $^1$H NMR spectra were

189 measured in deuterated chloroform.

190       Characteristic signals in the $^1$H NMR spectra of crude oil samples corresponding to

191 aromatic and aliphatic protons were found in the chemical shift regions 6.5–9.0 and 0.5–4.0

192 ppm, respectively (Fig. 2). As shown in Table 2, the aromaticity depends on the sample origin,

193 having values in the range of 2.10 % – 7.29 %. If compared with data summarized in Table 1,

194 $H_{ar}$ is well correlated with the asphaltene content, which is in agreement with the presence of

195 condensed aromatic rings in the asphaltene structure. On the other hand, no correlations were

196 observed as expected between the percentage of aliphatic chains calculated from the integral at

197 1.3 ppm ($I_{1,3ppm}$) and the asphaltene content, since aliphatic chains are present in all major crude

198 oil components.

199

200    **Fig. 2.** $^1$H NMR spectra of (a) Southwest Asia 1 and (b) North Africa 2 crude oil samples.

201    **Table 2.** Crude oil parameters calculated from $^1$H and DOSY NMR spectra.

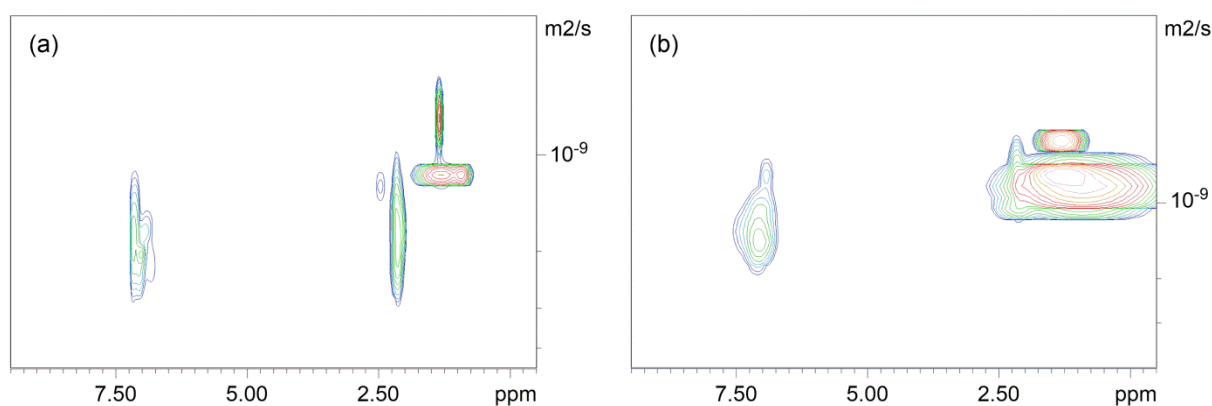| Sample No. | Designation | $H_{ar}$ / % [a] | $I_{1,3ppm}$ / % [b] | $d_{rel,0.9\ ppm}$ [c] | $d_{rel,1.3\ ppm}$ [d] |
|---|---|---|---|---|---|
| 1 | North Africa 1 | 5.70 | 56.75 | 0.52 | 0.51 |
| 2 | Southwest Asia 1 | 3.37 | 52.63 | 0.51 | 0.48 |
| 3 | Southwest Asia 2 | 4.24 | 56.49 | 0.55 | 0.53 |
| 4 | Southwest Asia 3 | 4.41 | 57.78 | 0.58 | 0.56 |
| 5 | Central Europe 1 | 7.29 | 61.25 | 0.53 | 0.50 |
| 6 | Eastern Europe 1 | 4.68 | 61.39 | 0.50 | 0.54 |
| 7 | Eastern Europe 2 | 3.69 | 59.59 | 0.54 | 0.54 |
| 8 | Central Europe 2 | 4.19 | 68.64 | 0.58 | 0.57 |
| 9 | West Africa | 4.54 | 44.83 | 0.50 | 0.52 |
| 10 | Southwest Asia 4 | 5.53 | 61.40 | 0.53 | 0.52 |
| 11 | Southwest Asia 5 | 5.66 | 53.16 | 0.55 | 0.54 |
| 12 | Eastern Europe 3 | 3.96 | 60.14 | 0.51 | 0.51 |
| 13 | NorthEast Asia | 4.84 | 59.48 | 0.57 | 0.56 |
| 14 | North Africa 2 | 2.10 | 58.92 | 0.80 | 0.74 |
| 15 | North Asia | 4.65 | 79.63 | 0.57 | 0.57 |
| 16 | Central Europe 3 | 3.49 | 69.94 | 0.51 | 0.51 |
| 17 | Central America | 2.82 | 63.58 | 0.77 | 0.52 |

202    [a]  aromaticity calculated as the difference between the sum of all signal integrals and those corresponding to
203        aliphatic protons
204    [b]  percentage of aliphatic chains obtained from the integral of the proton signal at 1.3 ppm
205    [c]  relative diffusivities calculated from the DOSY signal at 0.9 ppm
206    [d]  relative diffusivities calculated from the DOSY signal at 1.3 ppm

207     Further insight into the content and motional behavior of the crude oil components was

208     obtained from DOSY NMR experiments. Representative DOSY NMR spectra of crude oil

209     samples (Fig. 3) revealed differences in the shape and intensity of characteristic peaks. These

210     signals belong to species with different diffusion properties and can be used to distinguish

211     between the samples. Motional behavior of individual components is quantitatively described

212     by their translational diffusion coefficients ($D$). However, the accuracy and reproducibility of

213     the diffusion measurements is largely affected by experimental conditions. This impact can be

214     minimized by introducing the relative diffusivity, $d_{rel} = D_{sample}/D_{toluene}$. As shown for the signals

215     at 0.9 and 1.3 ppm in Table 2, as well as for other resonances in Table S1, only the components

216     of North Africa 2 and Central America samples exhibited considerably higher $d_{rel}$ than average.

217     On the other hand, diffusion properties of other crude oils were very similar to each other,

218     despite their different origin. Hence, in order to separate and classify all crude oils additional

219     information was extracted from DOSY NMR spectra by statistical analysis. For that purpose,

220     an approach was employed that combines MWA and ML methods described in the following
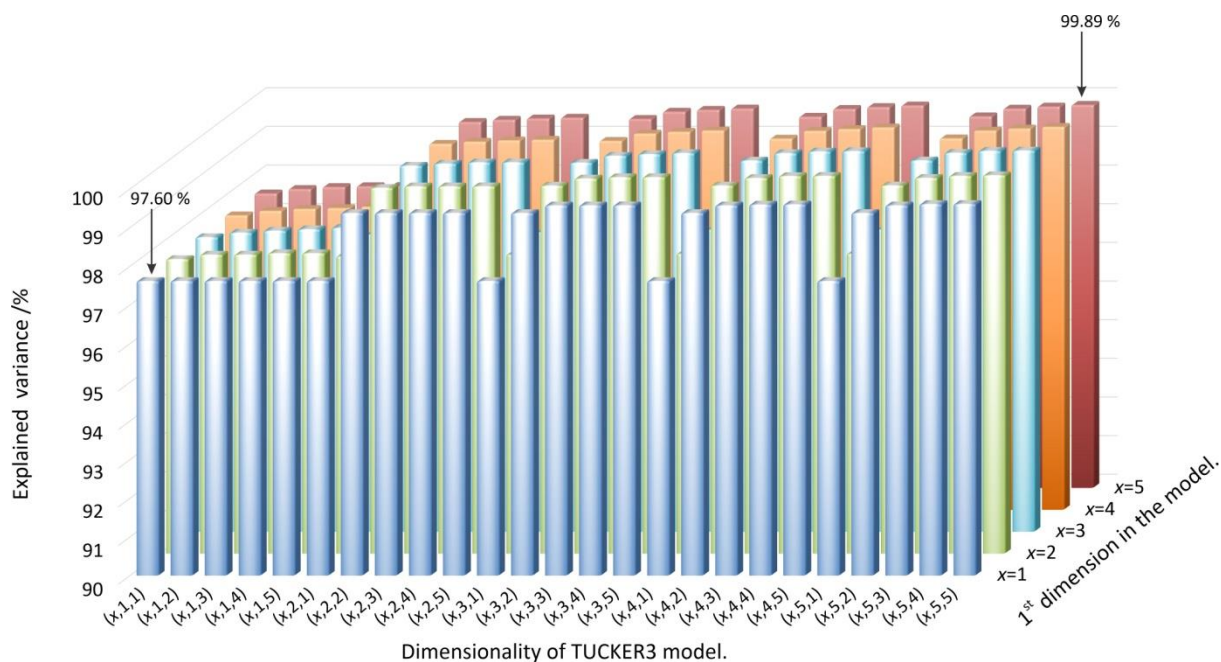
221     chapters.



**Fig. 3.** DOSY NMR spectra of (a) Southwest Asia 1 and (b) North Africa 2 crude oil samples.
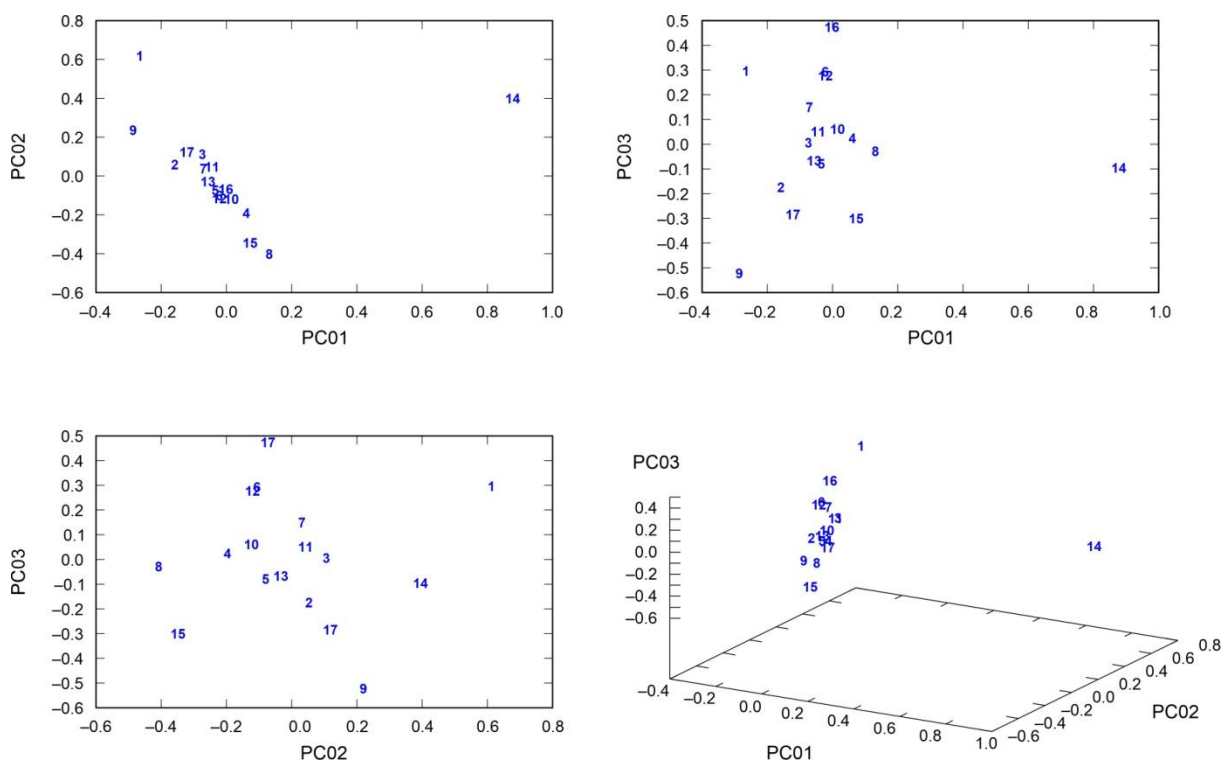
224     *3.3.    MWA*

225     Using the TUCKER3 decomposition model for a set of DOSY NMR spectra, a

226     *progressive* decomposition model search was performed starting from the model with

11

227 dimensions (1,1,1). This was the simplest decomposition model which already explained

228 97.60 % of the total variance (Fig. 4). The search passed through all possible models up to the

229 final tested decomposition model (5,5,5) that described 99.89 % of the total variance. Each

230 dimension was gradually increased by 1 giving the total number of generated models

231 $5 \times 5 \times 5 = 125$. Explained variances for all investigated models are presented in Fig. 4.



232

233 **Fig. 4.** Explained variance in TUCKER3 models in dependence of model dimensionality used

234 in decomposition of $3^{rd}$-order data tensor (DOSY NMR spectra).

235       Model (5,5,5) was chosen for further analysis, while the first three components from this

236 model were used for classification of samples, visualization and later regression. These three

237 components described 99.72 % of the total variance. Their loadings plots are presented in Fig.

238 5. This percentage of the total described variance is high enough to ensure that the most

239 important properties of the investigated systems relevant for the proper analysis were retained

240 within the model.

12

241

**Fig. 5.** Classification of the petroleum samples spanned in the space of the first three principal components for 3$^{rd}$-way loadings calculated by TUCKER3 decomposition.

In the reduced space of 3$^{rd}$-way loadings presented on Fig. 5, DOSY NMR spectra of the samples were represented as points (labeled as in Tables 1, 2 and S1). The distribution of all samples in this 3-dimensional space can be used for a classification of these samples based on their DOSY spectra. The sample **14** is highly distinguishable from the other samples, which is clear from the presented distribution. Moreover, from the variability among the samples one can see that the samples **6** and **12** are very similar. Investigation of 2-dimensional projections confirms that the same applies to the sample pairs (**10,11**) and (**5,13**).
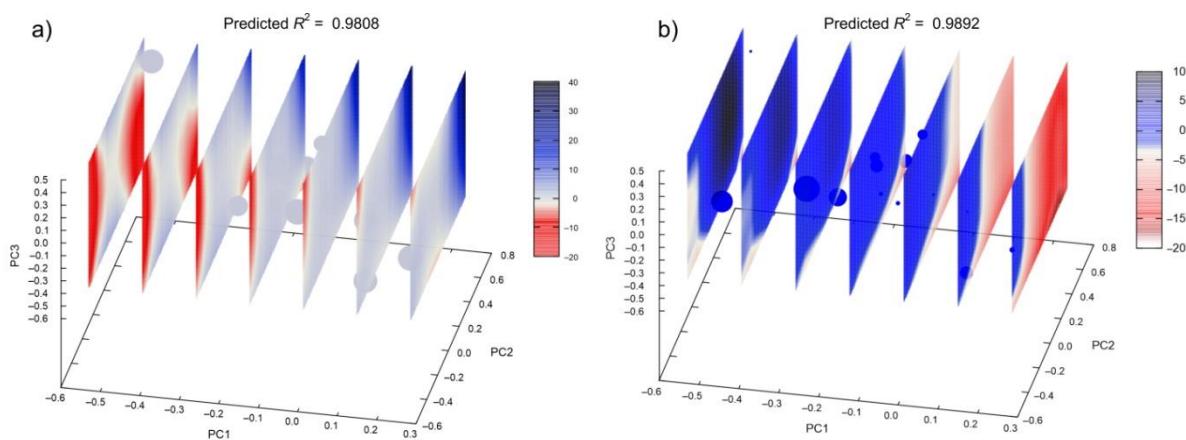
*3.4. Machine learning multivariate linear regression*

*3.4.1. Modeling stability with DOSY NMR spectra*

To establish a connection between measured stability data and DOSY NMR spectra, extensive ML procedure was utilized. Stability data $S_{asph}$, $S_{resin}$ and $S_{total}$ from Table 1 were regressed to the first three principal components in the 3-dimensional reduced space of DOSY

13

NMR spectra. In this way, each crude oil sample was represented by the point in the reduced three-dimensional space and stability was modeled using these three predictors as independent variables. The total number of generated different models for each dependent variable was 1 717 869 184 (models were built as linear combination of original variables, as well as their higher-order polynomial terms, up to the 4th order). The 4th order was shown to be sufficient for building excellent regression models. Models with polynomial terms up to the 3rd order had predicted $R^2$ of 0.89 (for all three measured stability values). It was therefore justified to push it up to the order of 4 judging the model quality on the basis of adjusted $R^2$, predicted $R^2$ and LOO-CV mean-squared-error ensuring that there was no overfitting.

The best determined models among all different 1 717 869 184 tested models had predicted $R^2$ to be bigger than 0.98. For the measured stability of asphaltenes $S_{asph}$, predicted $R^2$ was 0.9808 whereas for $S_{resin}$ predicted $R^2$ value was 0.9892 (Fig. 6). The quality of these models ensures that the crude oil stability in any similar crude oil sample can be predicted from the DOSY NMR spectra. This fact provides a broad range of possible applications using the DOSY NMR spectra for these or similar complex samples without the need for any additional chemical analyses. Properly predicting the stability of crude oils could *e.g.* directly reduce asphaltene remediation costs [29].



**Fig. 6.** The best multivariate regression model of the measured asphaltene and resin stability determined by machine learning: a) $S_{asph}$, and b) $S_{resin}$ in dependence on the 1st, 2nd and 3rd

276 principal component of the DOSY NMR spectra of crude oil samples obtained by MWA.

277 (Spheres represent points in 3D reduced space, and the planes are cuts of polynomial regression

278 model, for easier interpretation 4th-dimension is represented redundantly with the color and with

279 the size of the spheres.)

280 *3.4.2. Modeling stability with other measured properties*

281 Stability data was also regressed on 6 other measured properties: $H_{ar}$, $I_{1.3\,ppm}$, $w_{asph}$, $\rho_{API}$,

282 $d_{rel,0.9\,ppm}$ and $d_{rel,1.3\,ppm}$ (Tables 1 and 2). These measured properties were selected and their

283 selection was further confirmed by investigation of linear correlation matrix with measured

284 stability data where these properties showed some degree of linear correlation ($|R|>0.6$). In this

285 case the number of possible models with linear combination of terms up to the polynomial order

286 2 was 134 217 728. Using parallelized ML code [23], it was possible to test all these models

287 within one day and several excellent candidates were found. This search provided several

288 regression models with values of predicted $R^2$ higher than 0.99 for all three measured stability

289 parameters with the best ones having the following values of predicted $R^2$:

290 $R^2(S_{asph})=0.9998$, $R^2(S_{resin})=0.9997$ and $R^2(S_{total})=0.9999$.

291 These are particularly good values for such complex mixtures, proving that this new

292 model can accurately predict the crude oil stability and other important process parameters

293 relevant for petroleum industry.

294 Best determined models:

295 $S_{asph} = 1.73\text{E}{+}00\ +9.13\text{E}{-}01\times H_{ar}\ +1.99\text{E}{-}02\times I_{1.3\,ppm}\ -5.84\text{E}{-}02\times w_{asph}\ -3.39\text{E}{-}02\times \rho_{API}$
296 $-1.76\text{E}{+}01\times d_{rel,0.9\,ppm}\ -3.99\text{E}{-}02\times H_{ar}^2\ -1.21\text{E}{-}02\times H_{ar}\times \rho_{API}\ -1.40\text{E}{-}03\times I_{1.3\,ppm}^2$
297 $-1.11\text{E}{-}02\times I_{1.3\,ppm}\times w_{asph}\ +3.09\text{E}{-}01\times I_{1.3\,ppm}\times d_{rel,0.9\,ppm}\ +3.66\text{E}{-}02\times w_{asph}^2\ +1.06\text{E}{+}00\times$
298 $w_{asph}\times d_{rel,0.9\,ppm}\ +2.02\text{E}{-}01\times \rho_{API}\times d_{rel,1.3\,ppm}\ -5.03\text{E}{+}00\times d_{rel,0.9\,ppm}\times d_{rel,1.3\,ppm}$

299 $S_{resin} = -2.77\text{E}{+}01\ +1.52\text{E}{+}01\times w_{asph}\ +4.25\text{E}{+}00\times d_{rel,0.9\,ppm}\ +8.02\text{E}{+}01\times d_{rel,1.3\,ppm}$
300 $+7.83\text{E}{-}02\times H_{ar}^2\ -7.88\text{E}{-}01\times H_{ar}\times w_{asph}\ +3.75\text{E}{-}03\times I_{1.3\,ppm}^2\ +1.82\text{E}{-}02\times I_{1.3\,ppm}\times w_{asph}$
301 $-6.05\text{E}{-}03\times I_{1.3\,ppm}\times \rho_{API}\ -5.48\text{E}{-}01\times I_{1.3\,ppm}\times d_{rel,1.3\,ppm}\ -1.40\text{E}{-}01\times w_{asph}\times \rho_{API}$
302 $-1.49\text{E}{+}01\times w_{asph}\times d_{rel,0.9\,ppm}\ +1.85\text{E}{-}02\times \rho_{API}^2\ -1.27\text{E}{+}00\times \rho_{API}\times d_{rel,1.3\,ppm}$

303    $S_{total} = 1.78\text{E}+02 -9.21\text{E}+00\times w_{asph} -7.15\text{E}+00\times \rho_{API} -2.04\text{E}+02\times d_{rel,1.3\ ppm} +2.52\text{E}-02\times H_{ar}$

304    $\times I_{1.3\ ppm} -9.44\text{E}-02\times H_{ar} \times \rho_{API} +1.83\text{E}+00\times H_{ar} \times d_{rel,1.3\ ppm} -6.71\text{E}-04\times I_{1.3\ ppm}^2$

305    $-1.26\text{E}-02\times w_{asph} \times \rho_{API} +1.68\text{E}+01\times w_{asph} \times d_{rel,1.3\ ppm} +4.79\text{E}-03\times \rho_{API}^2 +1.32\text{E}+01\times \rho_{API} \times$

306    $d_{rel,0.9\ ppm} -5.59\text{E}+02\times d_{rel,0.9\ ppm}^2 +3.39\text{E}+02\times d_{rel,0.9\ ppm} \times d_{rel,1.3\ ppm}$

## 4. Conclusion

With the application of *multi-way* analysis using the TUCKER3 decomposition model for a set of DOSY NMR spectra, principal components were determined for the model (5,5,5). This decomposition model described 99.89% of the total variance. A classification of crude oil samples using the reduced space of the first 3 principal components was performed. Similar samples were identified and reduced space was further utilized for the regression of measured stabilities. Extensive machine learning multivariate linear regression was proven useful for modeling crude oil stability based on DOSY NMR spectra and other measured properties. For both cases, very good models were established, up to the 4[th] polynomial order in the first case and up to the 2[nd] polynomial order in the second one. This approach can serve as an excellent tool for predicting stability of complex petroleum samples and can be applied for other similar systems.

**Declaration of Competing Interest**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

325        **References**

326    [1]    Speight JG, El-Gendy NS. Chapter 1 - Petroleum composition and properties. In:

327            Speight JG, El-Gendy NS, editors. Introduction to petroleum biotechnology, Oxford:

328            Elsevier-Gulf Professional Publishing; 2018, p. 1–39. https://doi.org/10.1016/B978-0-

329            12-805151-1.00001-1.

330    [2]    Ashoori S, Sharifi M, Masoumi M, Salehi MM. The relationship between SARA

331            fractions and crude oil stability. Egypt J Pet 2017;26:209–13.

332            https://doi.org/10.1016/j.ejpe.2016.04.002.

333    [3]    Durand E, Clemancey M, Lancelin J-M, Verstraete J, Espinat D, Quoineaud, A-A.

334            Effect of chemical composition on asphaltenes aggregation. Energy Fuels

335            2010;24:1051–62. https://doi.org/10.1021/ef900599v.

336    [4]    Gawrys KL, Spiecker PM, Kilpatrick PK. The role of asphaltene solubility and

337            chemical composition on asphaltene aggregation. Petrol Sci Tech 2003;21:461–489.

338            https://doi.org/10.1081/LFT-120018533.

339    [5]    Evdokimov IN, Fesan AA, Losev AP. New answers to the optical interrogation of

340            asphaltenes: monomers and primary aggregates from steady-state fluorescence studies.

341            Energy Fuels 2016;30:4494–4503. https://doi.org/10.1021/acs.energyfuels.6b00027.

342    [6]    Rogel E, León O, Contreras E, Carbognani L, Torres G, Espidel J, et al. Assessment of

343            asphaltene stability in crude oils using conventional techniques. Energy Fuels

344            2003;17:1583–90. https://doi.org/10.1021/ef0301046.

345    [7]    Wang J, Buckley JS. Asphaltene stability in crude oil and aromatic solvents - the

346            influence of oil composition. Energy Fuels 2003;17:1445–51.

347            https://doi.org/10.1021/ef030030y.

348    [8]    Honse SO, Mansur CRE, Lucas EF. The influence of asphaltenes subfractions on the

349             stability of crude oil model emulsions. J Braz Chem Soc 2012;23:2204–10.

350             http://dx.doi.org/10.1590/S0103-50532013005000002.

351    [9]    Schermer WEM, Melein PMJ, van den Berg FGA. Simple techniques for evaluation of

352             crude oil compatibility. Pet Sci Technol 2004;22:1045–54. https://doi.org/10.1081/LFT-

353             120038695.

354   [10]   Vieira AP, Portela NA, Neto ÁC, Lacerda Jr. V, Romão W, Castro EVR, Filgueiras PR.

355             Determination of physicochemical properties of petroleum using [1]H NMR spectroscopy

356             combined with multivariate calibration. Fuel 2019;253:320–6.

357             https://doi.org/10.1016/j.fuel.2019.05.028.

358   [11]   Gao G, Cao J, Xu T, Zhang H, Zhang Y, Hu K. Nuclear magnetic resonance

359             spectroscopy of crude oil as proxies for oil source and thermal maturity based on [1]H

360             and [13]C spectra. Fuel 2020;271:117622. https://doi.org/10.1016/j.fuel.2020.117622.

361   [12]   Rakhmatullin I, Efimov S, Tyurin V, Gafurov M, Al-Muntaser A, Varfolomeev M et al.

362             Qualitative and quantitative analysis of heavy crude oil samples and their SARA

363             fractions with [13]C nuclear magnetic resonance. Processes 2020;8:995.

364             https://doi.org/10.3390/pr8080995.

365   [13]   Parlov Vuković J, Novak P, Jednačak T. NMR spectroscopy as a tool for studying

366             asphaltene composition. Croat Chem Acta 2019;92(3):323–29.

367             https://doi.org/10.5562/cca3543.

368   [14]   Durand E, Clemancey M, Lancelin J-M, Verstraete J, Espinat D, Quoineaud A-A.

369             Aggregation states of asphaltenes: Evidence of two chemical behaviors by [1]H diffusion-

370             ordered spectroscopy nuclear magnetic resonance. J Phys Chem C 2009;113:16266–

371             16276. https://doi.org/10.1021/jp901954b.

372 [15] Lisitza NV, Freed DE, Sen PN, Song Y-Q. Study of asphaltene nanoaggregation by

373 nuclear magnetic resonance (NMR). Energy Fuels 2009;23:1189–93.

374 https://doi.org/10.1021/ef800631a.

375 [16] Parlov Vuković J, Novak P, Jednačak T, Kveštak M, Kovačević D, Smrečki V, et al.

376 Magnetic field influence on asphaltene aggregation monitored by diffusion NMR

377 spectroscopy: Is aggregation reversible at high magnetic fields? J Disper Sci Technol

378 2020;41:179–87. https://doi.org/10.1080/01932691.2018.1561302

379 [17] Parlov Vuković J, Hrenar T, Novak P, Friedrich M, Plavec J. New multiway model for

380 identification of crude oil and asphaltene origin based on diffusion-ordered nuclear

381 magnetic resonance spectroscopy. Energy Fuels 2017;31:8095–8101.

382 https://doi.org/10.1021/acs.energyfuels.7b01358.

383 [18] Parlov Vuković J, Novak P, Plavec J, Friedrich M, Marinić Pajc Lj, Hrenar T. NMR and

384 chemometric characterization of vacuum residues and vacuum gas oils from crude oils

385 of different origin. Croat Chem Acta 2015;88:89–95. http://dx.doi.org/10.5562/cca2612.

386 [19] Standard Test Method for Determination of Asphaltenes (Heptane Insoluble) in Crude

387 Petroleum and Petroleum Products, ASTM D 6560-17.

388 [20] Standard Test Method for Determination of Intrinsic Stability of Asphaltene-Containing

389 Residues, Heavy Fuel Oils, and Crude Oils ($n$-Heptane Phase Separation; Optical

390 Detection), ASTM D 7157-18.

391 [21] Hrenar T, Primožič I, Fijan D, Majerić Elenkov M. Conformational analysis of spiro-

392 epoxides by principal component analysis of molecular dynamics trajectories. Phys

393 Chem Chem Phys 2017;19:31706–13. https://doi.org/10.1039/C7CP05600A.

394 [22] Tucker L. Some mathematical notes on three-mode factor analysis. Psychometrika

395 1966;31:279–311. https://doi.org/10.1007/BF02289464.

396   [23]  Hrenar T. *moonee,* Program for Manipulation and Analysis of Multi- and Univariate

397        Data. Revision 0.6827, Zagreb, Croatia, 2021.

398   [24]  Novak P, Kišić A, Hrenar T, Jednačak T, Miljanić S, Verbanec G. In-line reaction

399        monitoring of entacapone synthesis by Raman spectroscopy and multivariate analysis. J

400        Pharmaceut Biomed 2011;54:660–6. https://doi.org/10.1016/j.jpba.2010.10.012.

401   [25]  Jović O, Smolić T, Primožič I, Hrenar T. Spectroscopic and chemometric analysis of

402        binary and ternary edible oil mixtures: qualitative and quantitative study. Anal Chem

403        2016; 88:4516–4524. https://doi.org/10.1021/acs.analchem.6b00505.

404   [26]  Jović O, Smolić T, Jurišić Z, Meić Z, Hrenar T. Chemometric analysis of Croatian extra

405        virgin olive oils from central Dalmatia region. Croat Chem Acta 2013; 86:335–344.

406        http://dx.doi.org/10.5562/cca2377.

407   [27]  Parlov Vuković J, Telen S, Srića V, Novak P. The use of [13]C NMR spectroscopy and

408        comprehensive two-dimensional gas chromatography, GC×GC, for identification of

409        compounds involved in diesel fuel oxidative behavior. Croat Chem Acta 2011;84:537–

410        41. http://dx.doi.org/10.5562/cca1874.

411   [28]  IP 499/03 Standard Methods for Analysis and Testing of Petroleum and Related

412        Products and British Standard 2000 Parts, Methods IP 361 to 501. The Institute of

413        Petroleum London; 2003, 499.1.

414   [29]  Kraiwattanawong K, Fogler HS, Gharfeh SG, Singh P, Thomason WH, Chavadej S.

415        Thermodynamic solubility models to predict asphaltene instability in live crude oils.

416        Energy Fuels 2007;21:1248–55. https://doi.org/10.1021/ef060386k.