

# Analyst

Accepted Manuscript



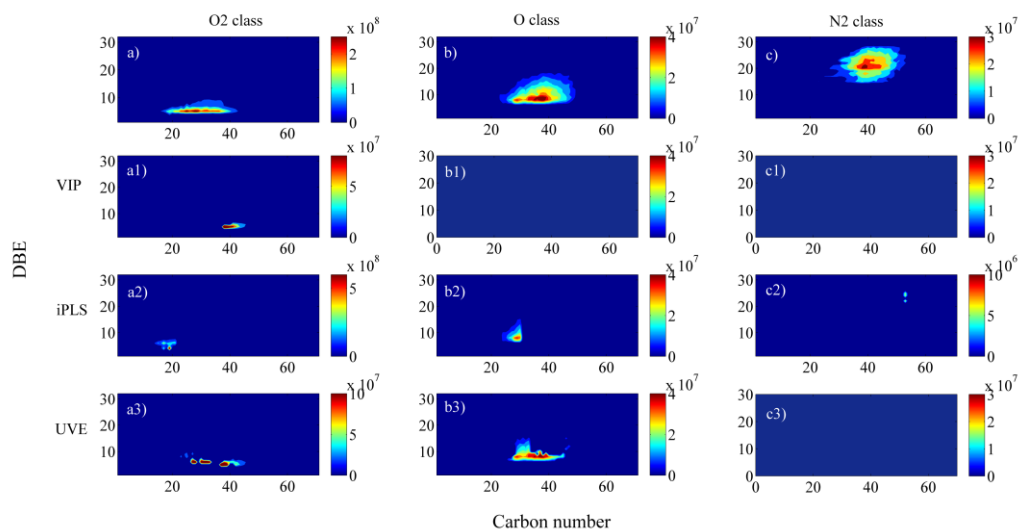
This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

*Accepted Manuscripts* are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

Fourier transform ion cyclotron resonance mass spectrometry (FT-ICR MS) was coupled to a Partial Least Squares (PLS) regression and variable selection methods to estimate the total acid number (TAN) of Brazilian crude oil samples



1  
2  
3 **Petroleomics by Electrospray Ionization FT-ICR Mass Spectrometry**  
4 **Coupled to Partial Least Squares with Variable Selection Methods:**  
5 **Prediction of the Total Acid Number of Crude Oils**  
6  
7  
8  
9

10 Luciana A. Terra<sup>a</sup>, Paulo R. Filgueiras<sup>a</sup>, Lílian V. Tose<sup>b</sup>, Wanderson Romão<sup>b,c</sup>, Douglas  
11 D. de Souza<sup>d</sup>, Eustáquio V. R. de Castro<sup>b</sup>, Mirela S. L. de Oliveira<sup>e</sup>, Júlio C. M. Dias<sup>e</sup>  
12 and Ronei J. Poppi<sup>a\*</sup>  
13  
14  
15

16  
17 <sup>a</sup> Institute of Chemistry, University of Campinas, Campinas, SP, Brazil.

18  
19 <sup>b</sup> Petroleomic and Forensic Chemistry Laboratory, Department of Chemistry, Federal  
20 University of Espírito Santo, 29075-910, Vitória, ES, Brazil.  
21  
22

23  
24 <sup>c</sup> Federal Institute of Education, Science and Technology of Espírito Santo, 29106-010,  
25 Vila Velha, ES, Brazil.  
26  
27

28  
29 <sup>d</sup> Institute of Physics “Gleb Wataghin”, University of Campinas, Campinas, SP, Brazil.

30  
31 <sup>e</sup> CENPES/PETROBRAS, Rio de Janeiro, RJ, Brazil.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 \* Corresponding author. Tel./fax: +55 19 35212134.  
54

55 E-mail address: ronei@iqm.unicamp.br (R. J. Poppi).  
56  
57  
58  
59  
60

**Abstract**

Negative-ion mode electrospray ionization, ESI(-), with Fourier transform ion cyclotron resonance mass spectrometry (FT-ICR MS) was coupled to a Partial Least Squares (PLS) regression and variable selection methods to estimate the total acid number (TAN) of Brazilian crude oil samples. Generally, ESI(-)-FT-ICR mass spectra presents a power of resolution of ca. 500.000 and a mass accuracy less than 1 ppm, producing a data matrix containing over 5700 variables per sample. These variables correspond to heteroatom-containing species detected as deprotonated molecules,  $[M-H]^-$  ions, which are identified primarily as naphthenic acids, phenols and carbazole analog species. The TAN values for all samples ranged from 0.06 to 3.61 mg of KOH g<sup>-1</sup>. To facilitate the spectral interpretation, three methods of variable selection were studied: variable importance in the projection (VIP), interval partial least squares (iPLS) and elimination of uninformative variables (UVE). The UVE method seems to be more appropriate for selecting important variables, reducing the dimension of the variables to 183 and producing a root mean square error of prediction of 0.32 mg of KOH g<sup>-1</sup>. By reducing the size of the data, it was possible to relate the selected variables with their corresponding molecular formulas, thus identifying the main chemical species responsible for the TAN values.

*Keywords:* ESI(-)-FT-ICR MS; Petroleomic; Total acid number; UVE-PLS.

## 1. Introduction

Petroleomics is defined as a field of petrol sciences able to elucidate the chemical composition of constituents present in crude oil due to its physical and chemical properties and reactivity. In general, crude oil is composed mainly of carbon (80 to 90%), hydrogen (10 to 15%), sulfur (up to 5%), oxygen (up to 4%), nitrogen (up to 2%) and traces of other elements (e.g., nickel and vanadium). The composition of the oil is classified in terms of the proportion of hydrocarbons and polar aromatic compounds.<sup>1</sup>

Oxygen-containing compound classes (naphthenic acids, O2 class; phenols, O1 class), are responsible for some undesirable properties, such as acidity (caused mainly by the presence of naphthenic acids), coloring, odors (phenols), the formation of emulsions and corrosion.<sup>2</sup> Among the main corrosive species, the naphthenic acids are evidenced in acidic crude oils, although they represent less than 3 wt %.<sup>3</sup> Naphthenic acids are defined as organic acids with the general formula  $R-(CH_2)_n-COOH$ , where R is a radical including one or more cyclopentane or cyclohexane rings (Figure 1). Due to naphthenic corrosivity effects and their biological marker role in geochemistry, many studies were focused on the identification of the naphthenic acids structures present in different crude oils.<sup>4</sup> This identification proved to be very difficult because naphthenic acids form complex mixtures. Some references mentioned that a single crude oil sample contains approximately 1500 different organic acids that are identified with molecular weights ranging from 200 to 700 Da.<sup>3,5</sup> More recent works on some crude oils identified naphthenic acid with a mass range of 115-1500 Da and a carbon content of C20 to C80.

The concentration of naphthenic acids in oils was one of the first tasks performed in the naphthenic corrosion studies. Currently, naphthenic concentrations are measured by titrating them with an alcoholic solution of potassium hydroxide (KOH), being

1  
2  
3 expressed by the total acid number (TAN) that represents the milligrams of KOH used  
4  
5 to neutralize all of the acidic species in 1 g of oil sample. Crude oils with TAN > 0.5 mg  
6  
7 of KOH/g may cause severe corrosion problems to refinery operations.<sup>6-8</sup> However, the  
8  
9 value of the TAN is not directly correlated to the corrosivity of naphthenic acids. The  
10  
11 TAN value depends upon the size and structure of the naphthenic acids and their  
12  
13 interaction with other compounds present in the crude oil (sulfites, carbon dioxide,  
14  
15 etc.).<sup>6</sup>  
16  
17

18  
19 Fourier transform ion cyclotron resonance mass spectrometry (FT-ICR MS)  
20  
21 offers the highest available mass resolution, mass resolving power and mass accuracy,  
22  
23 which enable the analysis of complex petroleum mixtures on a molecular level.<sup>9</sup> High-  
24  
25 resolution MS data have demonstrated that it is possible to discriminate many different  
26  
27 compounds<sup>10,11</sup> because of the different ionization efficiencies of the crude oil  
28  
29 constituents.<sup>12</sup> Accurate mass measurements<sup>13,14</sup> allow unambiguous elemental  
30  
31 composition ( $C_cH_hN_nO_oS_s$ ) assignment and DBE (double bond equivalents), facilitating  
32  
33 material classification by the heteroatom content and the degree of aromaticity.<sup>15,16</sup>  
34  
35 Naphthenic acids can be analyzed by negative-ion electrospray ionization, ESI(-),  
36  
37 coupled to FT-ICR MS, being detected in the form of a deprotonated molecule,  $[M -$   
38  
39  $H]^-$ .  
40  
41  
42

43  
44 The petroleomic MS characterization of crude oils has highlighted the  
45  
46 compositional trends to elucidate important crude oil properties. A fundamental goal of  
47  
48 petroleomics is to link such detailed crude oil compositions to its properties. To relate  
49  
50 the measured spectra to specific parameters, uni- and multi-variate calibration are often  
51  
52 used, which are especially useful with parameters that are difficult to measure  
53  
54 directly.<sup>17,18</sup> In addition, the use of mass spectrometry can reduce waste, minimizing the  
55  
56  
57  
58  
59  
60

1  
2  
3 consumption of raw materials and energy, thereby diminishing the environmental  
4  
5 impact.

6  
7 The aim of this work was the utilization of ESI(-)-FT-ICR MS technique in  
8  
9 conjunction to PLS regression to find a relationship between the mass spectra and the  
10  
11 TAN value. After that, by using variable selection methods, the dimension of the MS  
12  
13 data was reduced and, most important, the main chemical species responsible for the  
14  
15 relationship was identified. Then, it was possible to link the crude oil composition with  
16  
17 the TAN parameter that is the fundamental goal of petroleomics.  
18  
19

## 20 21 22 **2. Data Analysis**

### 23 24 25 **2.1. Partial Least Squares**

26  
27 In the partial least squares (PLS) approach, the matrix of instrumental responses  
28  
29 (X) is related to the vector of property of interest (y) by the liner relationship presented  
30  
31 in Equation 1:<sup>19</sup>

$$32 \quad \mathbf{y} = \mathbf{Xb} \quad (1)$$

33  
34 The arrays  $\mathbf{X}$  and  $\mathbf{y}$  are decomposed into latent variables, similar to the principal  
35  
36 component analysis (PCA)<sup>20</sup>, as represented by Equations 2 and 3,

$$37 \quad \mathbf{X} = \mathbf{TP}^t + \mathbf{E} \quad (2)$$

$$38 \quad \mathbf{y} = \mathbf{Tq}^t + \mathbf{f} \quad (3)$$

39  
40 where  $\mathbf{T}$  is the matrix of scores,  $\mathbf{P}^t$  and  $\mathbf{q}^t$  are the loadings, and  $\mathbf{E}$  and  $\mathbf{f}$  are the residues.  
41  
42 The scores  $\mathbf{T}$  are estimated from the weight coefficient  $\mathbf{W}$  that is obtained to minimize  
43  
44 the vector  $\mathbf{f}$  and to maximize the relationship between  $\mathbf{X}$  and  $\mathbf{y}$  given by Equation 1,

$$45 \quad \mathbf{T} = \mathbf{XW} \quad (4)$$

The regression coefficients used to relate  $\mathbf{X}$  and  $\mathbf{y}$ , are calculated according to Equation 5:

$$\mathbf{b} = \mathbf{W}(\mathbf{p}'\mathbf{W})^{-1}(\mathbf{T}'\mathbf{T})^{-1}\mathbf{T}'\mathbf{y} \quad (5)$$

## 2.2. Variables Selection Methods

The variable selection methods have been used to produce simpler, robust and interpretable models. In this paper, the variables selection methods tested were: uninformative variable elimination (UVE), variable importance in projection (VIP) and interval partial least squares (iPLS).

### 2.2.1. Uninformative variable elimination (UVE)

The UVE proposed by Centner<sup>21</sup> is a method of variable selection based on the reliability analysis of regression coefficients ( $\mathbf{b}$ ) that represent the contribution of each variable to the established model, which is calculated through a leave-one-out validation method. The principle of UVE is to add noise (artificial variables) to the matrix of instrumental responses and develop a PLS model for the data set containing the experimental and the artificial variables. The reliability criterion  $c_j$  (for each  $j$  variable) is calculated based on the ratio of the regression coefficient  $mean(b_j)$  and  $std(b_j)$ , which are the mean and the standard deviation, respectively, of the regression coefficients (Equation 6).

$$c_j = \frac{mean(b_j)}{std(b_j)}, \quad \text{for } j = 1, \dots, 2p$$

(6)



where  $p$  is the number of variables of the instrumental responses. To estimate a cutoff, i.e., the limit for the included and excluded variables, a criterion of informative or uninformative variable is formulated and presented in Equation 7:

$$cutoff = k * \max(abs(c_{noise})) \quad (7)$$

where  $k$  is an arbitrary value (normally 2),  $c_{noise}$  are the values for the artificial variables, and  $\max(abs(c_{noise}))$  is maximum absolute value of the reliability criterion.<sup>22</sup>

### 2.2.2. Variable importance in projection (VIP)

VIP is a combined measure of how much a variable contributes to describe the two datasets as PLS regression.<sup>23</sup> The idea is to accumulate the importance of each variable  $j$  being reflected by weight  $w_h$  from each latent variable. The VIP measure  $v_j$  is defined by the Equation 8,

$$v_j = \sqrt{p \sum_{h=1}^H [SS_h (w_{hj} / \|w_h\|^2)] / \sum_{h=1}^H (SS_h)} \quad (8)$$

where  $SS_h$  is the sum of squares described by the  $h$  latent variables. Hence, the  $v_j$  is a measure of the contribution of each variable according to the variance described by each PLS component, where  $(w_{hj} / \|w_h\|^2)$  represents the importance of the variable  $j$ . The  $v_j < 1$  (average of  $v$ ) condition indicates a non-important variable that could probably be removed.

### 2.2.3. Interval Partial Least Squares (iPLS)

1  
2  
3 The interval Partial Least Squares (iPLS) is a method introduced by Norgaard<sup>24</sup>,  
4  
5 in which the matrix of instrumental response is divided in a defined number of intervals,  
6  
7 and a PLS model is developed for each one of them. The sub-interval (or more than one  
8  
9 sub-interval) that presented the smallest cross-validated error was selected. The iPLS  
10  
11 can be an effective tool to determine the importance of different parts of a data set.  
12  
13

### 14 15 16 **3. Experimental**

#### 17 18 19 20 **3.1. Reagents and samples**

21  
22  
23  
24  
25 Anhydrous propanol, toluene and potassium hydroxide, each with a purity  
26  
27 higher than 99.5%, were purchase from Vetec, Brazil and used for the TAN  
28  
29 determinations. Ammonium hydroxide (NH<sub>4</sub>OH) and sodium trifluoroacetate (NaTFA)  
30  
31 were purchased from Sigma–Aldrich, USA and used for the ESI(-)-FT-ICR MS  
32  
33 measurements. In this study, thirty four crude oil samples from sedimentary basin of the  
34  
35 Brazilian coast were used.  
36  
37  
38  
39

#### 40 41 **3.2. TAN determination**

42  
43  
44  
45 TAN measurements were performed according to the standard ASTM method  
46  
47 (ASTM D664-09)<sup>25</sup> using a potentiometer (Metrohm Analytical Instruments and  
48  
49 Accessories, USA) with a reproducibility of pH  $\pm$  0.005 mg KOH g<sup>-1</sup>. In 250 mL  
50  
51 beakers, the crude oil samples (5  $\pm$  0.5 g) were weighed, dissolved in 125 mL of  
52  
53 water/anhydrous propan-2-ol/toluene (0.5/49.5/50 % in volume), and then titrated with a  
54  
55 0.1 mol L<sup>-1</sup> alcoholic KOH solution. Silver/Silver chloride (Ag/AgCl) reference  
56  
57  
58  
59  
60

1  
2  
3 electrode built into the same electrode body and a magnetic stirrer were used during the  
4  
5 titration. The results are expressed as milligrams of potassium hydroxide per gram of  
6  
7 sample required to titrate a sample in a solvent to a specified end point (see Table 1).  
8  
9 The TAN values ranged from 0.06 to 3.61 mg KOH g<sup>-1</sup>.  
10

### 11 12 13 14 **3.3. ESI(-)-FT-ICR MS measurements** 15

16  
17  
18 Petroleum samples were analyzed by ESI(-). Briefly, the crude oil samples were  
19  
20 diluted to  $\approx 1.2 \text{ mg mL}^{-1}$  in 50:50 (v/v) toluene/methanol (which contained 0.1% w/v of  
21  
22 NH<sub>4</sub>OH for ESI(-)). The resulting solution was directly infused at a flow rate of 5  $\mu\text{L}$   
23  
24 min<sup>-1</sup> into the ESI(-) source. The mass spectrometer (model 9.4 T Solarix, Bruker  
25  
26 Daltonics, Bremen, Germany) was set to operate over a mass range of  $m/z$  200-1300.  
27  
28 The ESI(-) source conditions were as follows: a nebulizer gas pressure of 1.0 bar, a  
29  
30 capillary voltage of 3.0-3.5 kV and a transfer capillary temperature of 250 °C. The ions  
31  
32 are accumulated in the hexapolar collision cell over a time period of 0.15-0.20 s,  
33  
34 followed by transport to the analyzer cell (ICR) through the multipole ion guide system  
35  
36 (another hexapole). Each spectrum was acquired by accumulating 200 scans of time-  
37  
38 domain transient signals in four mega-point time-domain data sets. The front and back  
39  
40 trapping voltages in the ICR cell were - 0.60 V and - 0.65 V, respectively, for ESI(-).  
41  
42 All mass spectra were externally calibrated using a NaTFA solution ( $m/z$  from 200-  
43  
44 1200) after they were internally recalibrated using a set of the most abundant  
45  
46 homologous alkylated compounds for each sample. A resolving power ( $m/\Delta m_{50\%} \approx 450$   
47  
48 000, in which  $\Delta m_{50\%}$  is the full peak width at half-maximum peak height) of  $m/z$  400  
49  
50 and a mass accuracy of < 1 ppm provided the unambiguous molecular formula  
51  
52 assignments for singly charged molecular ions. The mass spectra were acquired and  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 processed using a custom algorithm developed specifically for petroleum data  
4  
5 processing, Composer software (Sierra Analytics, Pasadena, CA, USA). The MS data  
6  
7 were processed, and the elemental compositions of the compounds were determined by  
8  
9 measuring the  $m/z$  values. To help visualize and interpret the MS data, a typical plot was  
10  
11 constructed, such as a heteroatomic-containing compounds diagram and a carbon  
12  
13 number *versus* the double bond equivalents (DBE) plot, where DBE is defined as the  
14  
15 number of rings plus the number of double bonds in a molecular structure. The  
16  
17 aromaticity of a petroleum component can be deduced directly from its DBE value  
18  
19 according to Equation 9,  
20  
21

$$22 \quad DBE = c - h/2 + n/2 + 1 \quad (9)$$

23  
24 where  $c$ ,  $h$  and  $n$  are the numbers of carbon, hydrogen and nitrogen atoms, respectively,  
25  
26 in the molecular formula.  
27  
28  
29  
30  
31

### 32 **3.4. Models development**

33  
34  
35  
36 After the mass attribution by the Composer software, the data set was  
37  
38 transformed to a data matrix in the software Matlab. The optimization of the PLS  
39  
40 models was performed using a k-fold cross-validation, and the best results were  
41  
42 obtained using auto-scaling as a preprocessing step. In this study, 25 samples for  
43  
44 calibration and nine samples for prediction were used. To facilitate the spectral  
45  
46 interpretation, three different variables selection methods were tested: UVE, VIP and  
47  
48 iPLS. The best method was chosen based on the error of prediction (RMSEP)  
49  
50 determined using Equation 10, the number of selected variables and the region of the  
51  
52 image corresponding to the original image obtained from the DBE vs. carbon number  
53  
54  
55  
56  
57  
58  
59  
60

plots. All of the PLS models were development in the software PLS Toolbox version 6.7 for Matlab.

$$\text{RMSEP} = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}} \quad (10)$$

where  $\hat{y}$  is the predicted TAN value by the PLS model and  $y$  is the reference TAN value given by ASTM D664-09.

## 4. Results and discussion

### 4.1. ESI(-) FT-ICR MS

Figures 2a-c displays the mass spectra of three typical Brazilian crude oils with different TANs: S1 (Figure 2a, TAN = 3.61 mg KOH g<sup>-1</sup>), S4 (Figure 2b, TAN = 1.90 mg KOH g<sup>-1</sup>) and S21 (Figure 2c, TAN = 0.17 mg KOH g<sup>-1</sup>). The ESI(-)-FT-ICR MS technique was able to identify 5703 variables per sample (the number of molecular formula assignments from the Composer software). In general, the ESI(-)-FT-ICR mass spectra have Gaussian profiles from  $m/z$  200 to 800 with an average molar mass distribution ( $M_w$ ) centered from 520 (sample S4, 2b) to 550 Da (sample 2c). Heteroatom-containing species were detected as deprotonated molecules,  $[M-H]^-$  ions, corresponding primarily to naphthenic acids, phenols and carbazole-analog species.

The naphthenic acid species concentration changes as a function of the TAN values. The enlarged area around  $m/z$  417 (Figures 2a-c) highlights the decrease of the intensity of the ion  $[C_{28}H_{50}O_2 - H]^-$ , which has an  $m/z$  of 417.3738 and DBE of 4; whereas the ions  $[C_{29}H_{38}O_2 - H]^-$  and  $[C_{30}H_{42}O - H]^-$  of  $m/z$  417.2799 and 417.3164, and DBEs of 11 and 10, respectively, increases in function of TAN.

1  
2  
3 Generally, petroleum samples have chemical compositions that differ  
4 significantly from one another. One way to display the similarities or differences  
5 between the signal patterns of crude oil samples is to build certain types of plots, such  
6 as plots of the relative abundances of different classes of compounds. In the present  
7 study, for the class profile diagrams shown in Figure 3, the relative amounts of each  
8 class were calculated by summing the abundances in one compound class and dividing  
9 by the total abundance of all species. The relative amounts of N, N<sub>2</sub>, NO, NO<sub>2</sub>, O and  
10 O<sub>2</sub> class compounds present in the samples S01, S04 and S021 are shown in Figure 3.  
11 For sample S01 of higher acidity (TAN = 3.61 mg KOH g<sup>-1</sup>), the O<sub>2</sub> class, composed  
12 primarily of naphthenic acids,<sup>26</sup> was the most abundant class. The N class (carbazole  
13 analog species) was the second most abundant, followed by the O and NO<sub>2</sub> classes  
14 (analogous to phenols and carbazoles with one carboxylic acid group or two hydroxylic  
15 groups). For sample S04 of intermediary TAN value, (TAN = 1.90 mg KOH g<sup>-1</sup>), an  
16 increase of the relative abundance of non-basic nitrogen-containing compound classes is  
17 observed (N, N<sub>2</sub>, NO and NO<sub>2</sub> classes), with the majority presence of phenol analogous  
18 compounds, O class, also evident for sample S21 of TAN = 0.17 mg KOH g<sup>-1</sup>. As  
19 consequence of the relationship observed between the TAN values and the relative  
20 abundance of the O<sub>2</sub> class, a plot of the relative abundance of the O<sub>2</sub>-containing acids  
21 species *versus* the TAN value was built for the 34 samples analyzed, as shown in Figure  
22 4, where a clear relation is observed similar to reported in literature.<sup>18</sup> Although these  
23 acids are the primary contributors to the overall acidity of the crude oils, other non-basic  
24 nitrogen compounds, phenols and non-polar sulfur components (not detected by ESI-  
25 FT-ICR MS), also contribute to the TAN value of the crude oils.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55

#### 56 4.2. Variable selection methods

57  
58  
59  
60

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

Figure 5a shows the original mass spectra and Figure 5b shows the spectra separated by classes, where the classes are organized in alphabetic order. This data reorganization was performed to facilitate the data analysis by chemometric methods and the results interpretation. Analyzing Figure 5, the need of data preprocessing is identified for the correction of the main differences between the intensity of the molecular ions. These differences were corrected by using auto-scaling as the preprocessing step, due to its performance in removing unwilling variation.

Figures 6, 7 and 8 show the procedure for variable selection by UVE, VIP and iPLS, respectively. For the UVE, a confidence level of 99% was defined, with six latent variables in the model, whereas for the VIP,  $VIP > 5$  was chosen as the threshold, with five latent variables. For iPLS, the spectra were split in intervals of size of 100 variables, and four intervals with five latent variables were selected to obtain the best results of the RMSEP. It is important to know that the O class is 3731 to 4346 and the O2 class is 4347 to 5000 compound numbers. Note that in all figures, the majority of variables selected were between these intervals, being the most important to be used to predict the TAN value in the rearranged spectrum (Figure 5b).

Subsequently, a new PLS model was developed with the different variable selection methods, and the results are presented in Table 2, including the number of latent variables and the number of variables selected along with the respective calibration and validation errors. The F-test is found to exhibit a similar accuracy for all of the methods because  $F_{exp} < F_{crit}$  with a confidence level of 95 %. Additionally, these errors were consistent with those obtained when all of the variables were used. However, now the model presents the advantage of the reduced number of variables generating more parsimonious models. Another important observation in Table 2 is that the number of

1  
2  
3 selected variables is similar between the VIP and UVE methods (the VIP and UVE  
4  
5 methods reduced the amount of 5703 variables to 159 and 183 variables, respectively).  
6

7  
8 The choice of the most representative variable selection method can be based on  
9  
10 the relationship of the selected variables with their respective molecular formula, thus  
11 identifying the species responsible for the TAN parameter (with the carbon number  
12 ranging from C<sub>15</sub> to C<sub>69</sub>, and the DBEs from 3 to 8). Therefore, it is possible to plot  
13  
14 DBE versus the carbon number using the selected variables, as presented in Figure 9,  
15  
16 for the O<sub>2</sub> (a) O (b) and N<sub>2</sub> (c) classes from the three chemometric methods used: VIP  
17  
18 (a<sub>1</sub>,b<sub>1</sub>,c<sub>1</sub>), iPLS (a<sub>2</sub>,b<sub>2</sub>,c<sub>2</sub>) and UVE (a<sub>3</sub>,b<sub>3</sub>,c<sub>3</sub>). For this plot, it was chosen a sample  
19  
20 with TAN of 1.18 mg KOH/g, an intermediated acidity value. For others samples, with  
21  
22 different TAN values, different figures could be obtained. It is possible to see, by the  
23  
24 observation of the images in Figure 9 that the region selected by UVE method for O<sub>2</sub>  
25  
26 and O classes corresponds to the correct range of carbon number in comparison with the  
27  
28 image with all variables. This observation demonstrates that the UVE method reduced  
29  
30 the number of variables with the same error of prediction of the other tested methods,  
31  
32 while retaining the chemical information; therefore, it can be considered to be the best  
33  
34 method for variables reduction.  
35  
36  
37  
38  
39

40  
41 The predicted TAN values obtained by UVE-PLS compared to the measured ones  
42  
43 by ASTM method are shown in Figure 10. Note the clear correlation between the value  
44  
45 of TAN obtained from ESI(-)-FT-ICR MS data and that obtained from conventional  
46  
47 method. The RMSEP of 0.32 mg KOH/g (for nine validation samples) is less than the  
48  
49 result obtained by the other studies<sup>18</sup> of 0.77 mg KOH/g (for 10 validation samples).  
50  
51  
52  
53

## 54 **5. Conclusions**

55  
56  
57  
58  
59  
60



1  
2  
3 ESI(-)FT-ICR MS coupled to multivariate calibration produces a powerful  
4 analytical tool to predict the values of TAN for different crude oils. This approach  
5 combines the reliable information supplied from ESI(-) FT-ICR MS data with the  
6 possibility of variables reduction and the multivariate calculation supplied by  
7 chemometric methods, such as PLS.  
8  
9

10  
11  
12  
13  
14 Generally, the ESI(-)-FT-ICR mass spectra produced a data matrix containing a  
15 total number of variables of 5703 for the different crude oils studied. Among the three  
16 variable selection methods studied (variable importance in the projection (VIP), interval  
17 partial least squares (iPLS) and elimination of uninformative variables (UVE)), the  
18 UVE was found to be more appropriate to select the important variables for the model,  
19 reducing the dimension of the variables from 5703 to 183. This result was confirmed  
20 from the DBE *versus* carbon number plot, where maximum distributions for the  
21 abundant species are in good agreement. The RMSEP obtained was 0.32 mg of KOH g<sup>-1</sup>,  
22 which is consistent with the error when using all variables. By reducing the size of the  
23 data, it was possible to relate the selected compounds with their corresponding  
24 molecular formula, thus identifying the relationship with the TAN of the oil.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37

38 In future works, other physical and chemical properties of crude oils will be  
39 predicted, such as density, viscosity, SARA method (saturate, aromatic, resin and  
40 asphaltene) and elemental analyses (percentage of carbon, hydrogen, nitrogen and  
41 oxygen), from non-polar compounds ionization, being then applied other ionization  
42 sources, such as laser desorption/ionization and atmospheric pressure photoionization.  
43  
44  
45  
46  
47  
48  
49  
50

## 51 Acknowledgments

52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 This research was generously funded by PETROBRAS/CENPES, CNPq,  
4  
5 CAPES, FAPESP and FINEP.  
6  
7  
8

### 9 10 **References**

- 11 1. B. P. Tissot and D. H. Welte, *Petroleum Formation and Occurrence*, Berlin, Springer-  
12 Verlag, 1984.  
13
- 14 2. J. G. Speight; *Handbook of Petroleum Analysis*, John Wiley and Sons, Inc.: New  
15 York, 2001.  
16
- 17 3. M. P. Barrow, J. V. Headley, , K. M. Peru, P. J. Derrick, *Energy Fuels*, 2009, **23**,  
18 2592-2599.  
19
- 20 4. R. Piehl, *Mater. Perform*, 1988, **27**, 37-43.  
21
- 22 5. L. B. Moura, R. F. Guimarães, H. F. G. Abreu, H. C. Miranda, S. S. M. Tavares,  
23 *Material Research*, 2012, **15**, 277-284.  
24
- 25 6. Q. Shi, S. Zhao, Z. Xu, K. H. Chung, Y. Zhang, C. Xu, *Energy Fuels*, 2010; **24**,  
26 4005–4011.  
27
- 28 7. E. Slavcheva, B. Shone, A. Turnbull, *Corros. J.*, 1999, **34**, 125-131.  
29
- 30 8. C. S. Hsu, G. J. Dechert, W. K. Robbins, E. K. Fukuda, *Energy Fuels*, 1999, **14**, 217-  
31 223.  
32
- 33 9. C. S. Hsu, C. L. Hendrickson, R. P. Rodgers, A. M. McKenna, A. G. Marshall, *J.*  
34 *Mass Spectrom.*, 2011, **46**, 337–343  
35
- 36 10. A. Gaspar, W. Schrader, *Rapid Commun Mass Spectrom.*, 2012, **26**, 1047-1052.  
37
- 38 11. A. M. McKenna, J. M. Purcell, R. P. Rodgers, A. G. Marshall, *Energy Fuels*, 2010,  
39 **24**, 2929-2938.  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 12. G. C. Klein, A. Angstrom, R. P. Rodgers, A. G. Marshall, *Energy Fuels*, 2006, **20**,  
4 668-672.  
5  
6  
7  
8 13. C. S. Hsu, C. L. Hendrickson, R. P. Rodgers, A. M. McKenna, A. G. Marshall, *J.*  
9 *Mass Spectrom.*, 2011, **46**, 337-343.  
10  
11  
12  
13 14. J. J. Savory, N. K. Kaiser, A. M. McKenna, F. Xian, G. T. Blakney, R. P. Rodgers,  
14 C. L. Hendrickson, A. G. Marshall, *Anal. Chem.*, 2011, **83**, 1732-1736.  
15  
16  
17  
18 15. R. L. Cunico, E. Y. Sheu, O. C. Mullins, *Pet. Sci. Technol.*, 2004, **22**, 787-798.  
19  
20  
21  
22 16. J. M. Purcell, I. Merdrignac, R. P. Rodgers, A. G. Marshall, T. Gauthier, I. Guibard,  
23 *Energy Fuels*, 2010, **24**, 2257-2265.  
24  
25  
26  
27 17. M. Benassi, A. Berisha, W. Romão, E. Babayev, A. Rompp, B. Spengler, *Rapid*  
28 *Commun Mass Spectrom.*, 2013, **27**, 825-834.  
29  
30  
31  
32 18. G. V. Boniek; V. A. Patrícia, F. C. R. Werickson, O. G. Alexandre; C. L. Rosana,  
33 Pereira, *Energy Fuels*, 2013, **27**, 1873-1880.  
34  
35  
36  
37 19. P. Geladi and B. R. Kowalski, *Anal. Chim. Acta*, 1986, **185**, 1-17.  
38  
39  
40  
41 20. S. Wold, K. Esbensen and P. Geladi, *Chemom. Intell. Lab. Syst.*, 1987, **2**, 37-52.  
42  
43  
44 21. V. Centner, D. Massart, O. E. Noord, S. Jong, B.M. Vandeginste, C. Sterna, *Anal.*  
45 *Chem.*, 1996, 68, 3851-3858.  
46  
47  
48  
49 22. Z. Xiaobo, Z. Jiewen, M. J.W. Povey, M. Holmes, M. Hanpin, *Anal. Chim. Acta.*,  
50 2010, **667**, 14-32.  
51  
52  
53  
54 23. C. M. Andersen and R. Bro, *Journal of chemometrics*, **24**, pp.728-737, 2010.  
55  
56  
57  
58  
59  
60

1  
2  
3 24. L. Norgaard, A. Saudland, J.Wagner, J. Nielsen, L.Munck, S. Engelsen, *Applied*  
4  
5 *Spectroscopy*, 2000, **54**, 413-419.  
6  
7

8 25. ASTM D664-09, Standard test Method for Acid Number of Petroleum Products by  
9  
10 Potentiometric Titration, ASTM International, West Conshohocken, PA, 2006.  
11

12 26. K. A. P. Colati, G. P. Dalmaschio, E. V. R. Castro, A. O. Gomes, B. G. Vaz, W.  
13  
14 Romão, *Fuel*, 2013, **108**, 647-655.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 1.** TAN values for the crude oil samples.

Samples	TAN (mg-KOH / g)	Samples	TAN (mg-KOH / g)
S 01	3.61	S 19	0.19
S 02	3.09	S 20	0.18
S 03	2.64	S 21	0.17
S 04	1.90	S 22	0.17
S 05	1.18	S 23	0.15
S 06	0.89	S 24	0.14
S 07	0.72	S 25	0.13
S 08	0.70	S 26	0.11
S 09	0.65	S 27	0.10
S 10	0.58	S 28	0.06
S 11	0.51	S 29	0.06
S 12	0.42	S 30	0.06
S 13	0.45	S 31	0.06
S 14	0.35	S 32	0.04
S 15	0.32	S 33	0.04
S 16	0.29	S 34	0.14
S 17	0.27		
S 18	0.24		

**Table 2.** RMSEP for the variable selection methods.

<b>Variable Selection Method</b>	<b>Number latent variables</b>	<b>Number select variables</b>	<b>RMSEC (mg-KOH/g)</b>	<b>RMSECV (mg-KOH/g)</b>	<b>RMSEP (mg-KOH/g)</b>	<b>*<math>R_{cal}^2</math></b>
None	6	5703	0.12	1.05	0.34	0.89
<i>VIP</i>	5	159	0.13	0.80	0.35	0.98
<i>iPLS</i>	6	400	0.34	0.56	0.38	0.87
<i>UVE</i>	6	183	0.20	0.79	0.32	0.93

\* Determination coefficient for calibration samples set.

## Captions for the Figures

**Figure 1.** Some chemical structures of the naphthenic acids present in crude oil.

**Figure 2.** ESI(-)-FT-ICR MS for the crude oil samples S01(a), S04 (b), e S21 (c). The insert shows that the intensity of ion  $[C_{28}H_{50}O_2 - H]^-$  of  $m/z$  417.3738 and DBE of four decreases, whereas the ions  $[C_{29}H_{38}O_2 - H]^-$  and  $[C_{30}H_{42}O - H]^-$  of  $m/z$  417.2799 and 417.3164, and DBEs of 11 and 10, respectively, increases in function of TAN.

**Figure 3.** Heteroatom-containing compound class distribution from the ESI(-)-FT-ICR MS data of the crude oil samples S01 (high acidity), S04 (intermediate acidity) and S21 (low acidity).

**Figure 4.** Percentage of the O2 class *versus* the TAN values for 34 crude oil samples.

**Figure 5.** ESI(-) FT-ICR mass spectra of 34 crude oils samples. Original mass spectrum (a) and mass spectrum rearranged by classes (b).

**Figure 6.** UVE results.  $t$  values for experimental (1-5703) and artificial random (5704-11406) variables (a) and selected variables (b). The cutoff level at 0.99 is indicated by the dashed line.

**Figure 7.** VIP variables of the PLS model.

1  
2  
3 **Figure 8.** Illustration of the results from iPLS. The columns denote the RMSECV  
4 obtained from each of the intervals. The horizontal line denotes the RMSECV of the  
5 full-spectrum model.  
6  
7  
8  
9

10  
11 **Figure 9.** DBE *versus* carbon number for O2 class (a), O class (b) and N2 class (c)  
12 using all variables and for variable selection methods: VIP (a1,b1,c1), iPLS (a2,b2,c2)  
13 and UVE (a3,b3,c3).  
14  
15  
16  
17  
18  
19

20 **Figure 10.** Plot of predicted by UVE-PLS *versus* the measurement by the ASTM  
21 method for the TAN values.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



Figure 1

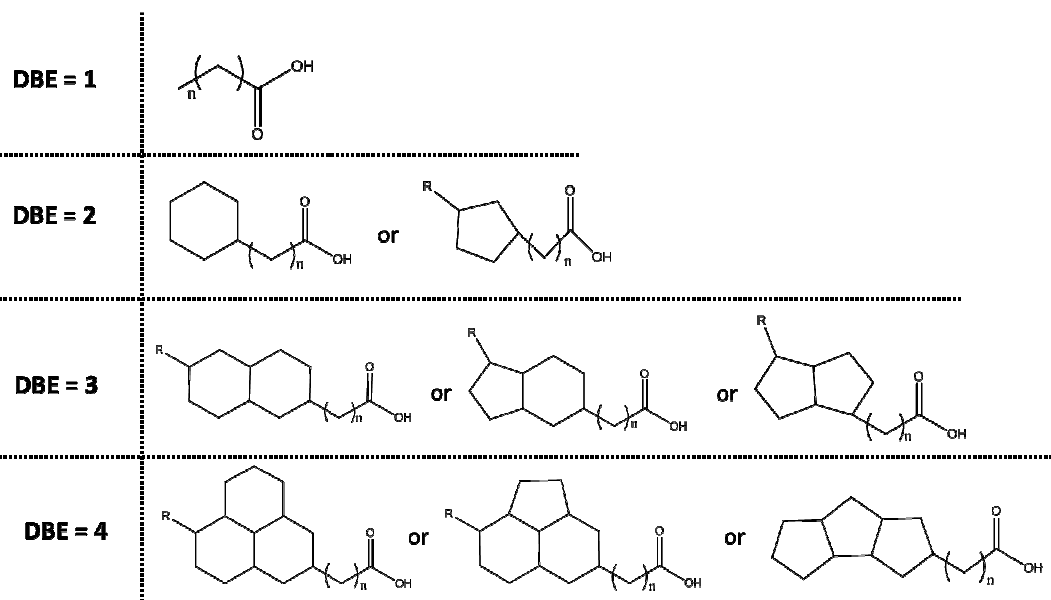


Figure 2

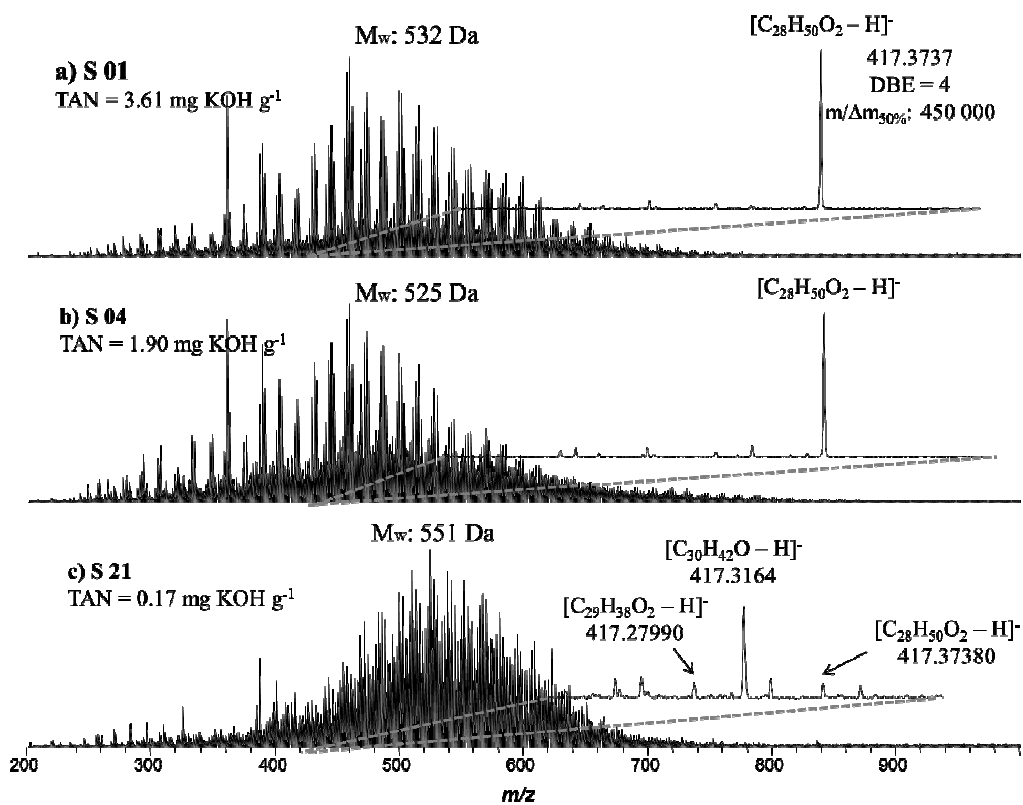


Figure 3

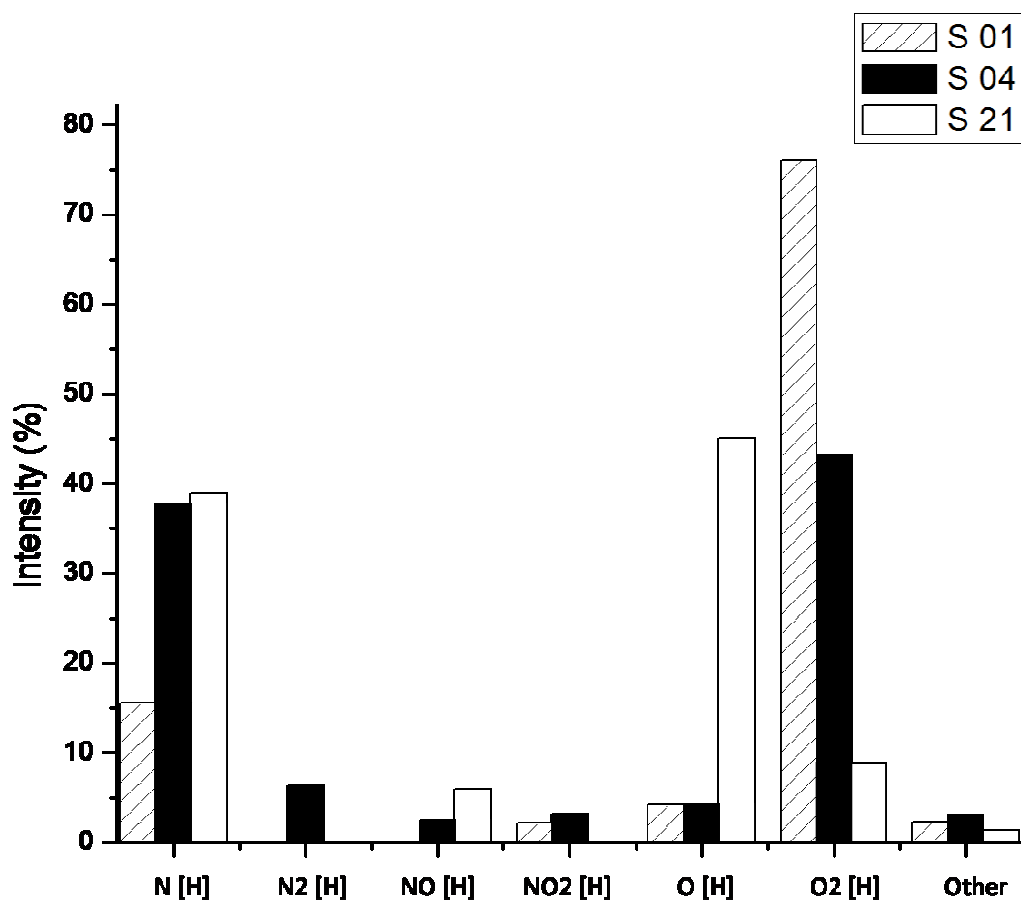


Figure 4

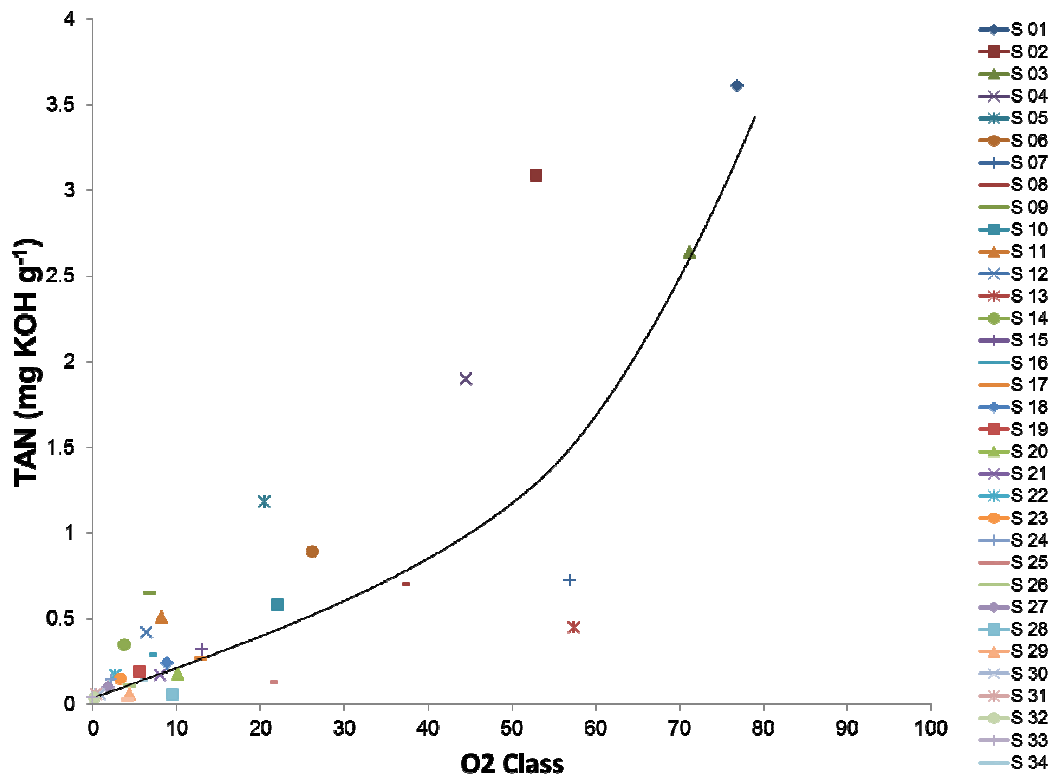


Figure 5

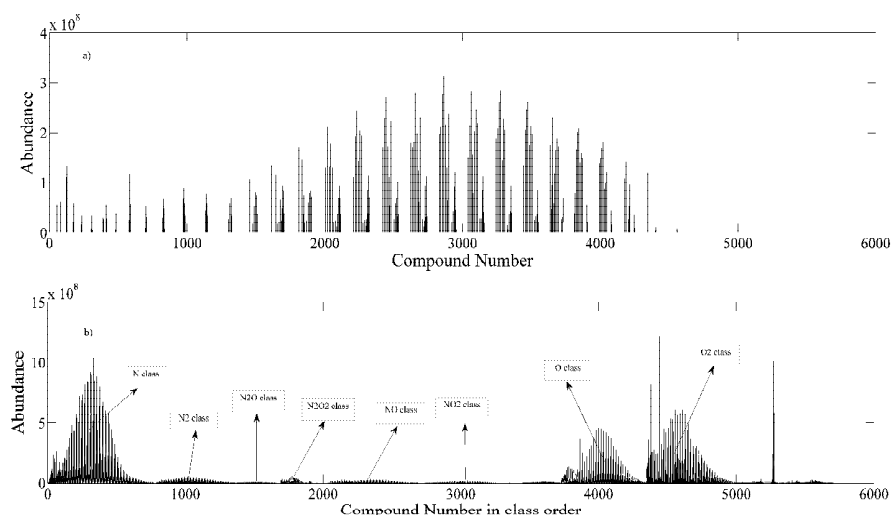


Figure 6

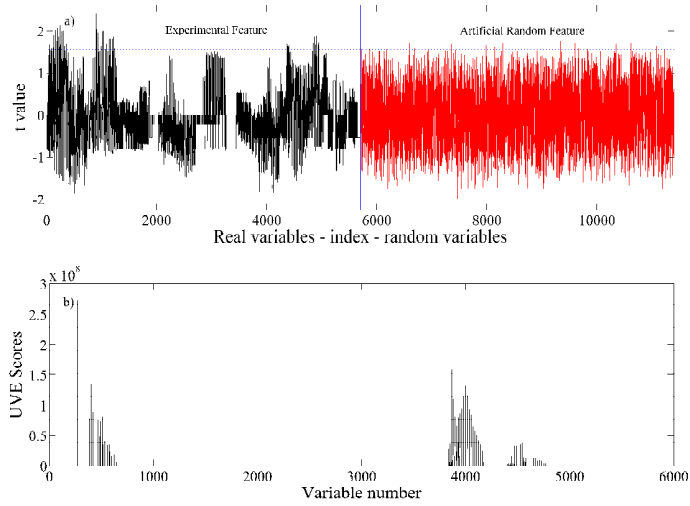
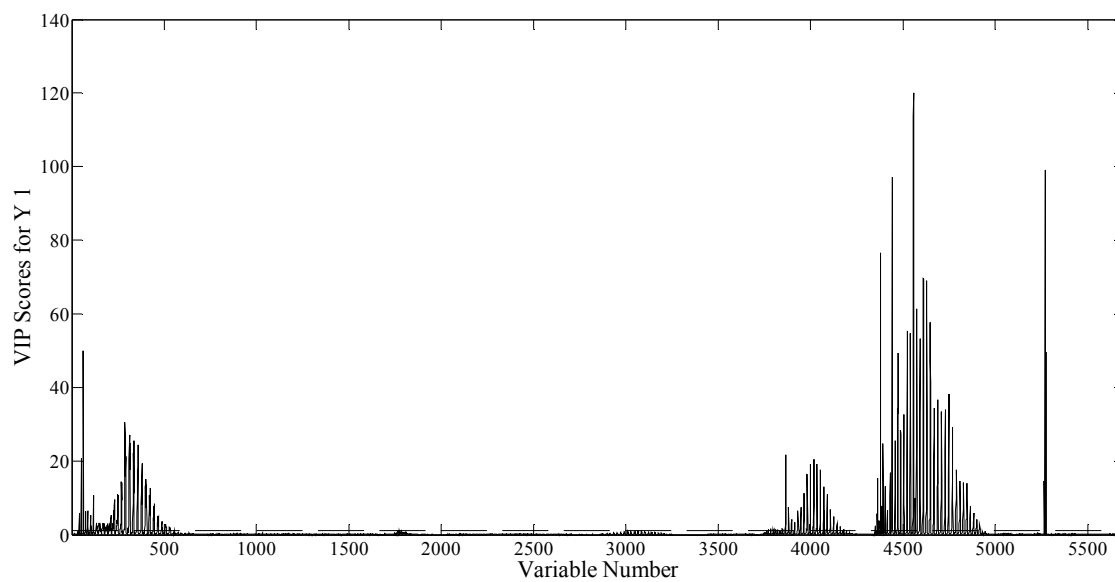
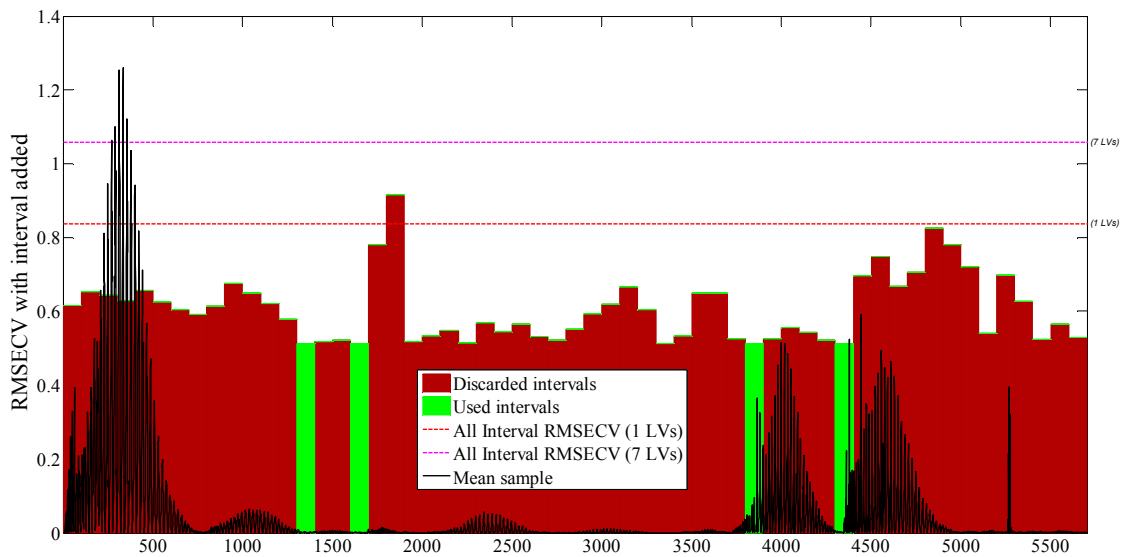


Figure 7



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

Figure 8



Analyst Accepted Manuscript



Figure 9

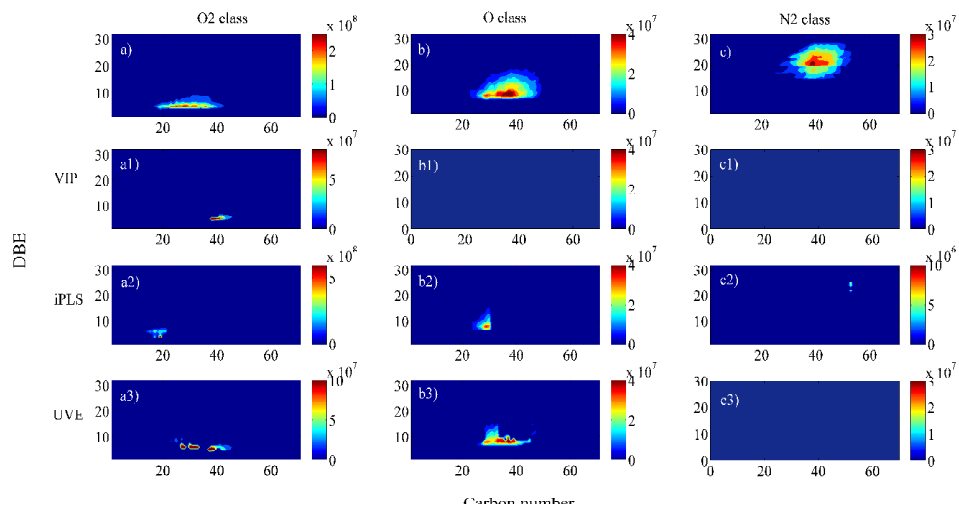


Figure 10

