

Significance of variables for discrimination - - applied to the search of organic ions in mass spectra measured on cometary particles

Kurt Varmuza^{1*}, Peter Filzmoser¹, Irene Hoffmann¹, Jan Walach¹,

Hervé Cottin², Nicolas Fray², Christelle Briois³, Johan Silén⁴, Oliver Stenzel⁵, Jochen Kassel⁵, Martin Hilchenbach⁵



¹Vienna University of Technology, Institute of Statistics and Mathematical Methods in Economics, Research Unit Computational Statistics, Vienna, Austria

²Université Paris-Est Créteil et Université Paris Diderot, Créteil, France; ³Université d'Orléans, Orléans, France

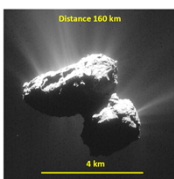
⁴Finnish Meteorological Institute, Helsinki, Finland; ⁵Max Planck Institute for Solar System Research (MPS), Göttingen, Germany



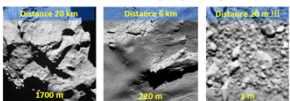
AIM. Recognition of mass spectral signals from CHNO (organic) compounds present at the surface of particles that have been collected near a comet (ESA mission Rosetta, TOF-SIMS instrument).

STRATEGY. Multivariate discrimination of spectra measured on a particle and on the background. Significance of variables indicates potentially relevant CHNO ions originating from the particle.

The Comet [1]



Name: **67P/Churyumov-Gerasimenko**
Orbit: 1.2 – 5.7 AU from sun, 6.4 years.
Rotation: 12.76 h. Density: 0.4-0.5 g/cm³.
Albedo: 5 % reflectance (very black).



The Mission (Rosetta) [2]

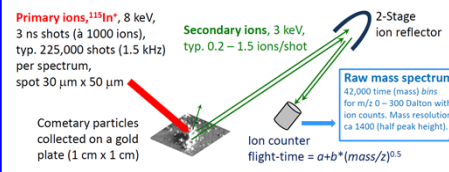


This mission of the European Space Agency (ESA) is named after the *Rosetta Stone* (196 BC), bearing an inscription that was the key to the decipherment of Egyptian hieroglyphs.



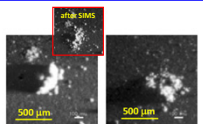
2 Mar 2004 Launch.
6 Aug 2014 Arrival (100 km).
Escorting the comet; typ. distance 10 – 200 km; 1.2 – 3.8 AU from sun.
12 Nov 2014 *Philae* landing.
30 Sep 2016 Controlled landing of Rosetta on comet (2 m/s), switch-off.

The Instrument (COSIMA) [3]



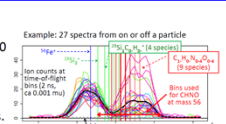
Samples

On metal targets (mostly Au) 1 cm x 1 cm more than 35,000 cometary particles were collected and imaged. Size 10 – 1000 μm. Fluffy material; some particles compact, others disrupted before collection [4].



Spectral Data

For about 200 particles more than 30,000 SIMS spectra were measured [5].
Preprocessing: Mass recalibration and rebinning; selection of mass bins (variables) for the considered CHNO ions.



Chemical Formulae of Ions

Considered mass range: 12 – 72 *

No. of formulae C_xH_yN_zO_w: 322

No. of mass bins (variables): 665 **

* Avoiding ions Si(CH₃)₃⁺, mass 73.047, from contamination PDMS (polydimethylsiloxane).

** Mass resolution 1000, ±1.5 σ (Gaussian peaks).

Formulae	No.
C H	50
C H N	112
C H O	85
C H N O	75
Total	322

Methods

Following methods have been used to characterize the importance of variables for a discrimination between comet particle spectra (class 1) and background spectra (class 2). A high **variable importance** indicates (via the mass corresponding to the variable) one or several ion species that may be prominent for the particle or for the background.

- **t-test** for comparing the class means (univariate; u-test gives almost identical results)
Data preprocessing: Normalization to constant sum of variables
Criterion: $LOGp = \text{sgn}[-\log(p)]$, p , probability of H₀; $\text{sgn} = +1$ if $\text{mean}_{\text{CLASS1}} > \text{mean}_{\text{CLASS2}}$, else $\text{sgn} = -1$
- **D-PLS**
Data preprocessing: Normalization to constant sum of variables
D-PLS: Optimum no. of PLS components estimated by repeated double cross validation (rdCV) [6]
Criterion: $bpls$ = standardized regression coefficients of discriminant variable (>0 indicates class 1)
- **Random Forest (RF)** classification
Data preprocessing: Normalization to constant sum of variables
RF: R library *randomForest*, function *randomForest*
Criterion: MDA (Mean Decreasing Accuracy, mean of 50 repetitions); a variable is important for class discrimination if the classification accuracy decreases considerably when the variable is eliminated.
- **Robust Pair-wise Log-Ratios (rPLR)**
Data preprocessing: Values $x_j = 0$ are replaced by uniformly distributed noise between 0.1 and 0.2 (median of all x is 33.4, maximum is 197).
rPLR: R library *robCompositions*, function *biomarker* [7]
Criterion: V^* is a newly developed statistics [7], derived from three variation matrices (for all objects, and for objects of the classes separately) containing robust estimations of the variances of all variable pairs (log-ratio transformed). The distribution of V^* is approximately standard normal, and quantile 0.975 is used to indicate the importance of the variables (biomarkers). The variable importance of a variable considers all other variables, and is not affected by the size effect, as ratios of variables are used.

Data

Positive second. ion spectra; measured at rectangular grid pos., $\Delta x = 15 \mu\text{m}$, $\Delta y = 10 - 30 \mu\text{m}$.
Mass range 12 – 72 Da
 $m = 665$ variables (mass bins)

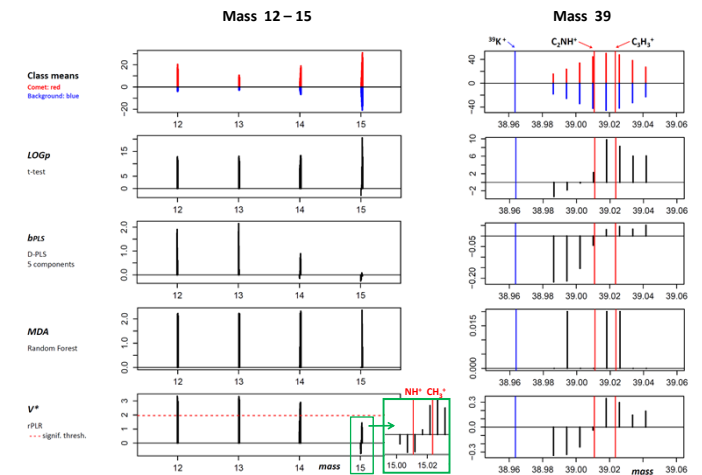
Spectra on particle: class 1
Particle *Sai*, target 3D1
Size ca 80 μm
Collected 2016-04-13/14
 $n_1 = 29$ spectra

Spectra backgr: class 2
Area *Nick*, target 3D0
Collected 2016-03-25/26
 $n_2 = 59$ spectra



Selected Results

Variable importance measure versus ion mass



Probably from the cometary material:

() not separable by mass; coded as *obvious* or as *guess*

C⁺; CH⁺; (CH₂⁺, N⁺); (CH₃⁺, NH⁺); C₂H₂⁺; (C₂H₃⁺, CNH⁺); C₃⁺; C₃H⁺; (C₃H₂⁺, C₂N⁺); (CH₂CN⁺, C₃H₄⁺); C₃H₃⁺; C₄⁺; ...

Probably from background:

Saturated or lowly unsaturated CH ions:
e. g., C₃H₅⁺, C₄H₇⁺, C₅H₇⁺

Summary

- Cometary particle surfaces contain CH(NO?) compounds.
- No distinct organic substance classes are evident from the data; a complex mixture of unsaturated organic compounds may be present.
- The results are consistent with the previously claimed presence of high molecular weight structures [8].
- The applied methods for characterizing the variable importance are complementary.
- Results from the univariate t-test are directly interpretable in terms of ion formulae. Results from the multivariate methods (D-PLS, RF, rPLR) characterize the variable importance within the potential influence of the other variables.

References

- http://www.esa.int/spaceinimages/Missions/Rosetta/class/image
- Schulz R., et al. (eds.): ROSETTA: ESA's mission to the origin of the solar system, Springer, New York (2009)
- Kassel J., et al.: Space Sci. Rev., **128**, 823 (2007)
- Langevin Y., et al.: Icarus, **271**, 76 (2016)
- Hilchenbach M., et al.: The Astrophysical Journal Letters, **816**, L32 (2016)
- Varmuza K., et al.: Chemom. Intell. Lab. Syst., **138**, 64 (2014)
- Walach J., et al.: submitted (2017)
- Fray N., et al.: Nature, **528**, 72 (2016)

This work is supported by the Austrian Science Fund (FWF), project P 26871 - N20.

Acknowledgments. COSIMA was built by a consortium led by the Max-Planck-Institut für Extraterrestrische Physik, Garching, Germany, in collaboration with the Laboratoire de Physique et Chimie de l'Environnement et de l'Espace, Orléans, France, the Institut d'Astrophysique Spatiale, CNRS/Université Paris Sud, Orsay, France, the Finnish Meteorological Institute, Helsinki, Finland, the Universität Wuppertal, Wuppertal, Germany, von Hoerner und Sulger GmbH, Schwetzingen, Germany, the Universität der Bundeswehr, Neuburg, Germany, the Institut für Physik, Forschungszentrum Seibersdorf, Seibersdorf, Austria, the Institut für Weltraumforschung, Österreichische Akademie der Wissenschaften, Graz, Austria and is led by the Max-Planck-Institut für Sonnensystemforschung, Göttingen, Germany. The support of the national funding agencies of Germany (DLR, grant 50QP1302), France (CNES), Austria, Finland and the ESA Technical Directorate is gratefully acknowledged. The authors thank the other members of the COSIMA team for their contributions.
J. W. thanks for support by the Austrian Science Fund (FWF) and the Czech Science Fund (GACR), project I1910-N26.