Atmospheric
Measurement
Techniques
Open Access
Discussions

EGU

1

1    **A Machine Learning Approach to Aerosol Classification for Single**

2    **Particle Mass Spectrometry**

3

4    **Christopoulos, Costa D.[1], Garimella, Sarvesh[1,2], Zawadowicz, Maria A.[1,3], Möhler,**

5    **Ottmar[4] and Cziczo, Daniel J.[1,5]**

6

7    [1] Department of Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of

8    Technology, Cambridge, MA, United States

9    [2] ACME AtronOmatic, LLC, Portland, OR, United States

10   [3] Atmospheric Sciences and Global Change Division, Pacific Northwest National

11   Laboratory, Richland, WA, United States

12   [4] Institute of Meteorology and Climate Research, Karlsruhe Institute of Technology,

13   Karlsruhe, Germany

14   [5] Department of Civil and Environmental Engineering, Massachusetts Institute of

15   Technology, Cambridge, MA, United States

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

2

**1**  **<u>Abstract</u>**

**2**  Compositional analysis of atmospheric and laboratory aerosols is often conducted via

**3**  single-particle mass spectrometry (SPMS), an *in situ* and real-time analytical technique

**4**  that produces mass spectra on a single particle basis. In this study, machine learning

**5**  classifiers are created using a dataset of SPMS spectra to automatically differentiate

**6**  particles on the basis of chemistry and size. Machine learning algorithms build a

**7**  predictive model from a training set for which the aerosol type associated with each mass

**8**  spectrum is known *a priori*. Classification models were also created to differentiate

**9**  aerosol within four broad categories: fertile soils, mineral/metallic particles, biological,

**10**  and all other aerosols. Differentiation was accomplished using ~40 positive and negative

**11**  spectral features. For the broad categorization, machine learning resulted in a

**12**  classification accuracy of ~93%. Classification of aerosols by specific type resulted in a

**13**  classification accuracy of ~87%. The 'trained' model was then applied to a 'blind'

**14**  mixture of aerosols which was known to to be a subset of the training set. Model

**15**  agreement was found on the presence of secondary organic aerosol, coated and uncoated

**16**  mineral dust and fertile soil.

**17**

**18**  **1. Introduction**

**19**  The interaction of atmospheric aerosols with clouds and radiation contributes to the

**20**  uncertainty in determinations of both anthropogenic and natural climate forcing [Boucher

**21**  et al., 2013; Lohmann and Feichter, 2005]. Aerosols directly affect atmospheric radiation

Atmospheric
Measurement
Techniques
Discussions

1    by scattering and absorption of radiation from both solar and terrestrial sources. The

2    radiative forcing from particulates in the atmosphere depends on optical properties that

3    vary significantly among different aerosol types [Lesins et al., 2002].   Aerosols also

4    indirectly affect climate via their role in the development and maintenance of clouds

5    [Vogelmann et al., 2012; Lubin et al., 2006]. Ultimately, the formation, appearance, and

6    lifetime of clouds are sensitive to aerosol properties like shape, chemistry, and

7    morphology [Lohmann and Feichter, 2008]. Characterization of aerosol properties,

8    therefore, plays a vital role in understanding weather and climate.

9        The chemical composition and size of aerosols has been analyzed on a single

10   particle basis *in situ* and in real-time using single particle mass spectrometry (SPMS;

11   Murphy [2007]). First developed ~2 decades ago, SPMS permits the analysis of aerosol

12   particles in the ~150 – 3000 nm size range, while differentiating internal and external

13   aerosol mixtures and characterizing both volatile (e.g. organics and sulfates) and

14   refractory (e.g. crystalline salts, elemental carbon and mineral dusts) particle components.

15   Particles are typically desorbed and ionized with a UV laser and resultant ions are

16   detected using time-of-flight mass spectrometry [Murphy, 2007]. A complete mass

17   spectrum of chemical components is normally produced from each analyzed aerosol

18   particle [Coe et al., 2006]. Despite almost universal detection of components found in

19   atmospheric aerosols, SPMS is not normally considered quantitative without specific

20   laboratory calibration [Cziczo et al., 2001].

21       Aerosols with different properties can appear similar in the context of SPMS. For

22   example, fly ash and mineral dust contain peaks corresponding to silicates, phosphates,

23   metals, and metal oxides despite different origins and emission sources [Zawadowicz et

Atmospheric
Measurement
Techniques
Open Access
Discussions
EGU

1    al., 2017]. This complicates analysis of aerosol populations because their properties need

2    to be well-defined in order to increase agreement between models and observations

3    [Niemand et al., 2012; Hoose and Möhler,  2012; Welti et al., 2009]. Even minor

4    compositional changes can be atmospherically important. As one example, mineral dusts

5    are known to be effective at nucleating ice clouds [Cziczo et al., 2013].  Particles in the

6    atmosphere undergo chemical and morphological changes as they age and eventually

7    contain material from several sources [Boucher et at. 2013]. Despite minor addition of

8    mass, aged mineral dust is less suitable for ice formation [Cziczo et al., 2013], but these

9    particles then act as cloud condensation nuclei and participate in warm cloud formation

10   [Andreae et al., 2008]. As a second example, ice nucleation in mixed-phase clouds has

11   been suggested to be predominantly influenced by feldspar, a single component among

12   the diverse mineralogy of atmospheric dust [Atkinson et al., 2013].

13        Here we show that supervised training and a rule-based probabilistic classification

14   of a decision tree ensemble can be used for differentiation of SPMS spectra. Various

15   clustering methods have been used to group aerosol types [Murphy et al., 2003; Gross et

16   al., 2008] but these algorithms are known to struggle with chemically-similar aerosols as

17   they do not incorporate known particle labels in the training process.  Such 'unsupervised'

18   clustering algorithms automatically group unlabeled data points on the basis of a

19   specified distance metric in feature space, in this case mass spectral signals. For the

20   purposes of setting broad aerosol categories, which are chemically similar and easily

21   separable in feature space, clustering is the simpler tool and the data easier to interpret.

22   For identifying new or potentially unexpected atmospheric aerosols, such properties are

23   desirable; however, the advantages of clustering greatly diminish when considering

Atmospheric
Measurement
Techniques
Discussions

1    similar particle types that overlap in feature space. Fertile soils, for instance, are often

2    grouped into a single category despite different sources and atmospheric histories.

3    Clustering algorithms should therefore be considered as a tool to use alongside

4    supervised classification. The latter may be used to further explore unique aerosol types

5    or verify manually labeled clusters with higher precision. Furthermore, the ensemble

6    approach presented here also produces variable rankings and probabilistic predictions that

7    assist in addressing measurement uncertainty.

8    In this study, we demonstrate the capabilities of machine learning to automatically

9    differentiate particles on the basis of chemistry and size. The resulting model can capture

10    minor compositional differences between aerosol mass spectra. By testing predictions

11    using an independent, or 'blind', dataset, we illustrate the feasibility of combining on-line

12    analysis techniques such as SPMS with machine learning to infer the behavior and origin

13    of aerosols in the laboratory and atmosphere.

14    **2. Methodologies**

15    **2.1 PALMS**

16    The Particle Analysis by Laser Mass Spectrometry (PALMS) instrument was

17    employed for these studies. PALMS has been described in detail previously [Cziczo et al.

18    2006]. Briefly, the instrument samples aerosol particles in the size range from ~200 to

19    ~3000 nm using an aerodynamic lens inlet into a differentially-pumped vacuum region.

20    Particle aerodynamic size is acquired by measuring particle transit time between two 532

21    nm continuous wave neodymium-doped yttrium aluminum garnet (Nd:YAG) laser beams.

22    A pulsed UV 193 nm excimer laser is used to desorb and ionize the particles and the

6

1    resulting ions are extracted using a unipolar time-of-flight mass spectrometer. The

2    resulting mass spectra correspond to single particles. The UV ionization extracts both

3    refractory and volatile components and allows analysis of all chemical components

4    present in atmospheric aerosol particles [Cziczo et al. 2013].

5

6    **2.2 Dataset**

7         A set of 'training data' was acquired by sampling atmospherically-relevant

8    aerosols. The majority of the dataset was acquired at the Karlsruhe Institute of

9    Technology (KIT) Aerosol Interactions and Dynamics in the Atmosphere (AIDA) facility

10   during the Fifth Ice Nucleation workshop — Part 1 (FIN01). The remainder were

11   acquired at our Aerosol and Cloud Laboratory at MIT. The FIN01 workshop was an

12   intercomparison effort of ~10 SPMS instruments, including PALMS. The training data

13   correspond to spectra of known particle types that were aerosolized into KIT's main

14   AIDA and a connected auxiliary chamber for sampling by PALMS and the other SPMSs

15   (Table 1). Hereafter we group both chambers with the name 'AIDA'. The number of

16   training spectra acquired varied by particle type, ranging from ~250 for secondary

17   organic aerosol (SOA) to ~1500 for potassium-rich feldspar ("K-feldspar"). In total,

18   ~50,000 spectra are considered with each spectrum containing 512 possible mass peaks

19   and an aerodynamic size. (Table 2). Additionally, the FIN01 workshop included a blind

20   sampling period, where AIDA was filled with 3 - 4 aerosol types known to be from the

21   training set (i.e., for which spectra had already been acquired) but (*a priori*) of unknown

22   size, specific types and at unknown concentrations.

1    Figure 1 illustrates a simple differentiation of particles using only two mass peaks

2    in one (negative) polarity. Mass peaks represent fractional ion abundance, measured as a

3    normalized total signal (ion current). In this example, the normalized areas of negative

4    mass peaks 24 ($C_2^-$) and 16 ($O^-$) are plotted. Distinct aerosol types are differentiated by

5    color with clusters forming in this two-dimensional space. Note that spectra of the same

6    aerosol type form distinct clusters (e.g. Arizona Test Dust, ATD), as do similar aerosol

7    classes (e.g., soil dusts). Co-plotted in Figure 1 are data from the blind experiment.

8    Distinct clusters of spectra from the blind experiment are noticeable and correlate with

9    known clusters.    Described in the next section, machine learning algorithms draw

10   "decision boundaries" that best separate different groups of data points based on set of

11   rules. Machine learning is not bound by the simplistic two dimensional space shown in

12   Figure 1 and instead uses all 512 mass peaks and aerodynamic size.

13   **2.3 Aerosol Classification**

14   A trained classification model maps a continuous input vector 'X' to a discreet

15   output value using a set of parameters 'learned' from the data. Figure 2 illustrates the

16   mapping of a mass spectrum to vector space. In contrast to traditional, hard-coded, rule-

17   based classification methods, machine learning determines parameters that partition the

18   data set. To form X, mass spectra are converted to dimensional vectors normalized to the

19   total ion current (i.e., the total of all mass peaks sum to 1 in each spectrum). The elements

20   of the vectorized mass spectrum, termed 'features', hold information about the ionization

21   efficiency and relative abundance of chemical species in each aerosol and serve as the

22   variables for the machine learning model.

1    Machine learning is conducted in two phases: training and testing. During training,

2    a model is constructed and iteratively updated based on data (i.e., mass spectra) from the

3    training set. For this work, the set of known aerosol types sampled by PALMS was

4    converted to dimensional vectors. These data form the basis set for defining each aerosol

5    type. An ensemble of decision trees was used to generate predictions of aerosol type. A

6    single decision tree is a statistical decision model that performs classification based on a

7    series of comparisons relating a variable $X_i$ (in this case a normalized mass peak in X) to

8    a learned threshold value [Breiman, 2001]. Represented as an algorithmic tree, a binary

9    decision tree consists of a hierarchy of nodes where each node connects via branches to

10   two other nodes deeper in the tree. At each node, one of the two branches is taken based

11   on whether a normalized peak $X_i$ is greater or less than a threshold value. Each branch

12   leads to another node where a different test is performed. After a series of tests, one at

13   each node, a class is assigned to a given sample; these are the so-called 'leaves'. Figure 2

14   illustrates the classification model for a single decision tree.

15   Each test in the tree narrows the set of reachable output leaves and thus the

16   sample space of possible aerosol labels. After $h$ tests in this study, where $h$ ranges from

17   10 to 3000, the set of reachable leaves and possible labels is 1 and the decision tree

18   outputs a prediction. Because PALMS is unipolar – either a positive or negative mass

19   spectrum is produced – simultaneous generation of positive and negative spectra on a

20   particle-by-particle basis is not possible.  Two separate classification models, one for

21   each polarity, were therefore generated to classify aerosols. These are hereafter referred

22   to as the 'positive' and 'negative classification algorithms'.

23

Atmospheric
Measurement
Techniques
Discussions

Open Access

1    **2.4 Decision Tree Ensembles**

2        An ensemble consists of a collection of classifiers where each independently

3    labels a spectrum vector X. To make a final prediction of aerosol type, decision trees

4    within an ensemble 'vote' on a classification label. Each vote has equal weight and the

5    spectrum is assigned to the majority choice. Each tree within an ensemble is

6    independently grown on a subset of the training data so that a commonly voted label

7    implies a higher certainty. Adding members to an ensemble increases the robustness of a

8    classification model by providing alternative hypotheses and is therefore preferable to

9    single classifiers.

10        Before an ensemble method is implemented for classification, trees are

11    independently grown during training.  A total of $k$ trees, with $k = 1000$, were grown using

12    a bootstrap sample from the training set. In bootstrap sampling, each tree sees an

13    independent sample set of equal size drawn from the full training set by sampling spectra

14    with replacement. On average, each tree is built with ~63% of the data. The unsampled

15    data, known as 'out-of-bag' observations, provide a means to assess classification error

16    for each tree during the training process.

17        Given a bootstrap sample, a binary decision tree is grown by sequentially creating

18    tests that maximize the separation between classes in parameter space.  A test is created

19    by defining a comparison that minimizes the information entropy of a possible split, thus

20    minimizing the randomness of prediction labels [Breiman, 1996]. To generate variability

21    in the model, a best split is chosen among a random set of possible splits at each node on

22    the basis of entropy [Breiman, 2001]. After iteratively defining thresholds for each new

1    node, the tree grows in size until a series of tests ending at some node $S_q$ uniquely

2    characterizes an aerosol as a particle type. A leaf is then appended to node $S_q$ with the

3    corresponding label. In classification mode, an aerosol spectrum that passes the same tree

4    will undergo the same series of tests and will end in the same leaf, thus being labeled in

5    the same way. For the purposes of this study, each tree had ~3,300 nodes.

6    **2.5 Dimensionality Reduction and Chemical Feature Selection**

7        Dimensionality reduction is the process of representing data with fewer variables

8    than initially present in the dataset, in this case less than the original 512 mass peaks and

9    aerodynamic size. In addition to facilitating data visualization, reducing computation time

10    and limiting overfitting [Mjolsnes, 2001], dimensionality reduction, in the context of

11    aerosol mass spectra, also indicates the most important chemical makers for

12    differentiation. Feature ranking was algorithmically determined by comparing the

13    performance of trees before and after removing information about peak $X_i$. The method is

14    that the values of variable $X_i$ is permuted for tree $k$ in the out-of-bag set so that the

15    variable is irrelevant to the final label. The change in misclassification before and after

16    the permutation is calculated and then repeated for all trees so that a variable ranking is

17    obtained [Breimann, 2001]. Table 2 rows ranks mass peaks (features) by polarity in

18    importance using this method. The columns at left list feature rankings (i.e., most to least

19    important for correct classification) for the entire set of aerosol types. The columns at

20    right list rankings when aerosol types are grouped into the broad, chemically similar,

21    categories. A final ranking was determined by sequentially adding variables and

22    observing classification performance response. All variables preceding two e-foldings in

23    classification error were maintained in the final model. Both the specific aerosol type and

Atmospheric
Measurement
Techniques
Discussions

1    broad aerosol category models were retrained using this subset of the initial variables,

2    listed in Table 2.

3        It is noteworthy that while most of the features are logical differentiators of the

4    aerosol types investigated in FIN01 there were also surprises. One example is $59^+$

5    (cobalt), determined to be one of the most important features for differentiation. Further

6    investigation determined this material was a contaminant from dry powder dispersion

7    equipment used on some samples. This serves to illustrate the lack of *a priori* judgment

8    by the algorithm and an unintended benefit of machine learning process (i.e.,

9    contamination identification).

10       **3. Results**

11   **3.1 Confusion Matrices and Probabilistic Model Performance**

12       A confusion matrix captures misclassification tendencies by pair-wise matching

13   the model prediction with the true aerosol type or broad category [Powers, 2007].

14   Confusion matrices represent model predictions as columns $i$ and true aerosol type of

15   category as rows $j$, where class names are mapped to integers $i, j \in \{1, 2, \dots, y\}$. In this

16   study, matrices have been normalized along each column to show the fraction of aerosols

17   labeled as $j$ that actually belong to $i$ (Figures 3 and 4). For aerosol classification, these

18   matrices can also be interpreted as similarity measures between particle types. Since the

19   basis of decision tree classification is mathematical separation of physical quantities,

20   misclassifications result from similarity in mass peaks and their ion abundance between

21   aerosol types. This is most easily visualized as overlapping clusters in the simple two

22   dimensional space in Figure 1.

1        Because the size of the set is large (~22,300), the general classification behavior

2    can be quantified in term of conditional probability. If $\hat{Y}_i$ is the set of predicted aerosol

3    spectra with aerosol label $i$ and $Y_j$ is the corresponding set of true spectrum-label pairs for

4    label $j$, then the conditional probability of assigning an aerosol to label $i$ given a predicted

5    label $j$ is given by:

6                          $$p(i \mid j) = \frac{|Y_j \cap \hat{Y}_i|}{|Y_j|} \qquad\qquad (1)$$

7        C is the raw confusion matrix of spectrum counts and $p(i \mid j)$ is the conditional

8    probability distribution over all true aerosol labels $i$, conditioned on some model-

9    generated label $j$. To obtain matrix P, which encodes $p(i \mid j)$ for all possible labeling,

10    columns of C are normalized with respect to the total aerosol counts for each label with

11    Eq. 1.

12        Model performance for each aerosol is summarized in the diagonal elements of P,

13    which represent the fraction of aerosol in column j labeled correctly. The classification

14    accuracy ($a$) is given by averaging diagonal elements of P. A perfect classification model

15    produces the identity matrix, as all data points are classified correctly 100% of the time.

16    For example, in the positive confusion matrix, SOA and Agar growth medium are

17    correctly labeled in the test set 100% of the time. Barring element truncation, all columns

18    of P add to 1.

19        Figures 3 and 4 display confusion matrices as heat maps for the full set of particle

20    labels and broad grouped particle categories, respectively. Broad categories are

21    delineated by bold horizontal and vertical lines in Figure 3 as fertile soil (Argentinian,

22    Chinese, Ethiopian, Moroccan and two German soils), pure mineral dust and metallic

23    particles (ATD, illite NX, fly ash, Na-feldspar, K-feldspar), biological (Agar growth

1   medium, *P. syringae* bacteria, cellulose, Snomax, and hazelnut pollen), and other (K-

2   feldspar with sulfuric acid (SA) and SOA coatings, soot, and SOA) particles. Some

3   model confusion exists between fertile soils and coated/uncoated feldspars which can be

4   explained since soils are mineral dust mixed with organic and other materials.

5       Positive mass spectra appear to hold more information with respect to

6   differentiating aerosols than negative. Label-wise classification accuracy for the negative

7   algorithm ranges from 3-5% lower. A large part of this performance discrepancy is due to

8   greater ability of positive spectra to differentiate coated particles within the 'other'

9   category.

10      In addition to quantifying misclassification tendencies between classes, the

11  confusion matrix can be redefined to show confusion for aerosols within broad categories

12  themselves. Intraclass misclassification analysis is accomplished by considering smaller

13  portions of C and using the same probabilistic assumptions highlighted for the full

14  confusion matrices to form modified probability distributions. The full confusion matrix

15  is partitioned into submatrices representing confusion in a specific aerosol category and

16  renormalized with respect to matrix columns. L is the subset of particle labels of a

17  broader set of aerosols.  Integrating the full conditional probability distribution over

18  labels that are impossible to observe gives the probability distribution over members of

19  L:

$$P_l(i,j) = p(i \in L \mid j \in L) = \frac{C(i \in L, j \in L)}{\sum_{i' \in L} C(i', j \in L)} \qquad (2)$$

20      For example, to determine $P_l(i \mid j)$ for fertile soils, a submatrix is formed by

21  collecting spectral counts in the first 6 rows and columns of the full confusion matrix

22  (Figure 3). Column normalization is then applied to derive a probability distribution over

Atmospheric
Measurement
Techniques
Discussions

Open Access

1    labels in the fertile soil category, conditioned on the aerosol actually being a fertile soil.

2    This analysis is repeated over all categories in both models. Finally, the relative

3    performance of both models is isolated and considered with respect to each specific

4    aerosol category.

5         The precision score [Powers, 2007] captures the classification behavior for some

6    subset of aerosol L by averaging fractions of correctly classified aerosols for labels

7    within that category:

8              $$\text{Precision Score(L)} = \frac{1}{|L|}\sum_{i=j}^{|L|} P(\, i \,\in L, j \in L) \qquad (3)$$

9

10        When applied to $P_l$, the precision score captures classification performance on a

11    population with only aerosol labels contained in L. The algorithm is expected to correctly

12    label an aerosol in such a population with a probability equal to the precision score. The

13    precision score is valuable when using the classification model as a particle screener,

14    producing probability distributions over a subset of aerosol labels of interest. The

15    confusion characteristics are shown in Table 3 for each category in terms of the precision

16    score and the mean and standard deviation of misclassification within each category.

17    Although both models perform similarly for biological spectra, discrepancies of 2-5%

18    appear in the remaining categories. For regimes consisting of only mineral/metallic or

19    other particles, the positive algorithm shows intraclass performance advantages in terms

20    of the precision score, but most notably in terms of fewer mislabeling of mineral/metallic

21    particles.   The largest precision discrepancy is observed for fertile soils, where the

22    positive ion algorithm has a 5% advantage in precision with approximately half the false

23    labeling rate.

24

Atmospheric
Measurement
Techniques
Open Access
Discussions
EGU

1   **3.2 Characterization of Blind Data**

2   As part of the FIN01 workshop, it was known that 3 - 4 aerosol types from Table

3   1 were aerosolized into the ADIA chamber but at unknown size and relative

4   concentration. PALMS, one member of the blind intercomparison effort, collected

5   ~25,000 spectra. After data analysis, the aerosol types and relative abundances were

6   provided to each group (Figure 5, top center).

7   The presence or absence of particle types in the blind set was initially diagnosed

8   by choosing particles predicted at or above the 1% level. We note here that this step was

9   based on the knowledge that (1) a distinct set of particles would be placed in the chamber

10  and (2) particles present at or below the 1% level were most likely contamination. We

11  further note that this step is unique to a blind study and would not be applicable to the

12  atmosphere. Normalized confusion matrices were redefined for the aerosols in the

13  population (i.e., those above the 1% level), which forms the labels of set L in Eq. 2.

14  Finally, particle counts are re-computed by reassigning particle labels based on the

15  modified confusion matrix. For each particle label $j$, a fraction $n' = P(i\,|\,j)$ of particles

16  labeled as $j$ are reassigned to $i$. This probabilistic correction accounts for aerosol

17  mislabeling tendencies observed during testing, producing statistics that better represent

18  the underlying aerosol population. The expected fraction of particles belonging to label $i$

19  (denoted $\hat{n}i$) is given by:

20  $$< \hat{n}i > = \frac{<n_i>}{|n|} = \frac{1}{|n|} \sum_j P(i\,|\,j)\,|n_j| \qquad (4)$$

21  where $n$ is a set containing all blind spectra and $n_j$ is the set of particles labeled as $j$.

22  Figure 5 illustrates the results after this step, where the bottom charts show

23  corrected fractional percentages for each aerosol category. Because SOA was nearly

Atmospheric
Measurement
Techniques
Discussions

Open Access

1   always labeled correctly (Figure 3), the remaining aerosols are considered separately

2   using the full set of candidate aerosol labels. Both positive and negative models arrived at

3   similar results, with inconsistencies primarily associated with the presence of trace fertile

4   soils and mineral dust / fly ash particles. The positive algorithm identifies ~2-4% of the

5   AIDA population as each Argentinean soil, German soil, ATD, and cellulose whereas the

6   frequency of these aerosols was too low to consider in the negative. Alternatively, the

7   negative model estimates Na-Feldspar at ~8% of the total population, a label not

8   identified by the positive algorithm. This discrepancy can be explained by the 1%

9   selection criterion for aerosols present in the population. Fertile soils, ATD, and cellulose

10  frequently accumulate error along rows in the full positive confusion matrix, indicating

11  frequent confusion with other categories (Figure 3). Furthermore, with the observed

12  misclassification rates ranging ~1-4%, it is expected that these aerosol labels are false

13  positives. The negative model offers an alternative hypothesis, suggesting these

14  miscellaneous aerosols are Na-feldspar. Since there is significant model agreement on the

15  percentages of SOA, K-Feldspar, and coated feldspars, this part of the blind mixture

16  population (~90%) can be characterized with most certainty. For the disputed aerosol

17  labels, more credence is lent to the negative classification algorithm on the basis of

18  improved precision for fertile soils.

19      The aerosols reported in the blind mixture were soot, mineral dust, and SOA. This

20  mineral component was not defined and may have been either a specific mineral or soil

21  dust. The soot aerosols were below the cutoff diameter for PALMS; they were therefore

22  not detected or identified by the algorithms. Similarly, particles with diameters greater

23  than ~1000 nm are detected with increasingly large inefficiency which likely leads to

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

1    undercounting of mineral dust [Cziczo et al., 2006]. Both algorithms robustly labeled

2    SOA with large agreement, consistent with the 100% accuracy observed in the test set.

3            SOA coated mineral dust was identified as a particle type. This material was not

4    directly input to AIDA but the report is most likely correct, due to coagulation within the

5    AIDA chamber during the course of the blind experiment. The training data set did not

6    contain coagulated SOA and mineral dust but did include SOA-coated K-Feldspar, which

7    explains the identification.

8            While both models identified a variety of fertile soils, and not a single type, these

9    results are largely consistent with the known uncertainties highlighted by the confusion

10   matrices discussed previously. Given the presence of any single mineral dust, some

11   confusion with fertile soils, SA coated Feldspar, and Na-Feldspar is expected (Figure 3).

12   Moreover, as discussed previously [Gallavardin et al., 2008], AIDA backgrounds are not

13   completely particle-free. During the FIN01 study, contamination particles from previous

14   test aerosol were frequently observed as background and they could also be the origin of

15   some low-concentration particles matching fertile soil chemistry.


16   **4. Conclusions and Future Work**

17           The machine learning approach described here allows for differentiation of

18   aerosols within a SPMS dataset, augmenting existing tools and reducing the need for a

19   qualitative comparison between mass spectra. This study lays out a framework for

20   training and implementing an ensemble classification model and interpreting results in

21   the context of laboratory and atmospheric aerosol populations. Across a representative

22   sample of possible aerosol types, the behavior of each algorithm predictably allows users

1    to infer the presence or absence of specific aerosols and quantify aerosol abundance.

2    Machine learning is automated and the output of the model must then be informed by

3    human knowledge of aerosol chemistry. Machine learning should therefore be considered

4    as an additional tool to interpret mass spectra to better distinguish aerosols with unique

5    properties in terms of atmospheric chemistry, biogenic cycles, and population health.

6    The ensemble decision tree classification framework described here may be

7    generalized to any instrument, or set of instruments, capable of collecting physical and

8    chemical information that distinguishes particles. Although the method described here is

9    applied to a stand-alone SPMS and tested with a set of 'blind' data, ancillary laboratory

10   or field data can be integrated to expand the data set. The success of these algorithms is

11   data-dependent, where better performance is expected for instruments that provide more,

12   and more quantitative, analysis of the aerosol properties. Although the algorithms

13   implemented in this study were primarily used to categorize SOA, mineral dust, fertile

14   soil and biological aerosols, these models can adopt an arbitrary large set of aerosol data.

15   **Acknowledgements**

21

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

# 1 References

2 Andreae, M. & Rosenfeld, D.: Aerosol–cloud–precipitation interactions. Part 1. The

3 nature and sources of cloud-active aerosols, Earth-Sci. Rev., 89, 13-41,

4 doi:10.1016/j.earscirev.2008.03.001, 2008.

5 Atkinson, J., Murray, B., Woodhouse, M., Whale, T., Baustian, K., & Carslaw, K.,

6 Dobbie, S., O'Sullivan, D., and Malkin, T. L: The importance of feldspar for ice

7 nucleation by mineral dust in mixed-phase clouds, Nature, 498, 355-358,

8 doi:10.1038/nature12278, 2013.

9 Boucher, O., Randall, D., Artaxo, P., Bretherton, C., Feingold, G., Forster, P., Kerminen,

10 V.-M. , Kondo, Y., Liao, H., Lohmann, U., Rasch, P., Satheesh, S.K., Sherwood,

11 S., Stevens B., and Zhang, X. Y.,: Clouds and Aerosols, Climate Change 2013:

12 The Physical Science Basis. Contribution of Working Group I to the Fifth

13 Assessment Report of the Intergovernmental Panel on Climate Change, 5, 571-

14 657, 2013.

15 Breiman L.: Bagging Predictors. Machine Learning, 24, 123-140, 1996.

16 Breiman L.: Random Forests. Machine Learning, 45, 5-32, 2001.

17 Coe, H., Allan, J. D.: In Analytical Techniques for Atmospheric Measurement; Heard, D.

18 E., Ed., Blackwell Publishing, 265–311, 2006.

19 Cziczo, D., Thomson, D., Thompson, T., DeMott, P., and Murphy, D.: Particle analysis

20 by laser mass spectrometry (PALMS) studies of ice nuclei and other low number

21 density particles, Int. J. Mass. Spectrom., 258, 21-29, 2006.

Atmospheric
Measurement
Techniques
Discussions

1  Cziczo, D. J., Froyd, K., Hoose, C., Jensen, E., Diao, M., Zondlo, M., Smith, J. B.,

2      Twohy, C. H., and Murphy, D. M.: Clarifying the Dominant Sources and

3      Mechanisms of Cirrus Cloud Formation, Science, 340, 1320-1324,

4      doi:10.1126/science.1234145, 2013.

5  Cziczo, D. J., Thomson, D. S., and Murphy, D. M.: Ablation, flux, and atmospheric

6      implications of meteors inferred from stratospheric aerosol, Science, 291 (5509),

7      1772–1775, 2001.

8  Gallavardin, S., Lohmann, U., and Cziczo, D.: Analysis and differentiation of mineral

9      dust by single particle laser mass spectrometry, Int. J. Mass. Spectrom., 274,

10     56-63, doi:10.1016/j.ijms.2008.04.031, 2008.

11  Gallavardin, S. J., Froyd, K. D., Lohmann, U., Möhler, O., Murphy, D. M., Cziczo, D. J.:

12     Single Particle Laser Mass Spectrometry Applied to Differential Ice Nucleation

13     Experiments at the AIDA Chamber, Aerosol Sci. Tech., 42, 773-791, doi:

14     10.1080/02786820802339538, 2008.

15  Garimella, S., Wolf, M. J., Christopoulos, C. D., Zawadowicz, M. A., and Cziczo, D. J.:

16     Measuring the cloud formation potential of fly ash particle, Atmos. Chem. Phys.

17     (in prep)

18  Gross, D., Atlas, R., Rzeszotarski, J., Turetsky, E., Christensen, J., Benzaid, S., Olson, J.,

19     Smith, T., Steinberg, L., and Sulman, J.: Environmental chemistry through

20     intelligent atmospheric data analysis, Environ. Modell. Softw., 25,

21     760-769, 2008.

22  Henning, S., Ziese, M., Kiselev, A., Saathoff, H., Möhler, O., Mentel, T. F.,

23     Buchholz, A., Spindler, C., Michaud, V., Monier, M., Sellegri, K. and

Atmospheric
Measurement
Techniques
Open Access
EGU
Discussions

1    Stratmann, F.: Hygroscopic growth and droplet activation of soot

2    particles: uncoated, succinct or sulfuric acid coated, Atmos. Chem. Phys.,

3    12(10), 4525–4537, doi:10.5194/acp-12-4525-2012, 2012.

4

5    Hoose, C. and Möhler, O.: Heterogeneous ice nucleation on atmospheric aerosols: a

6    review of results from laboratory experiments, Atmos. Chem. Phys., 12, 9817-

7    9858, doi:10.5194/acpd-12-12531-2012, 2012.

8    Hiranuma, N., Augustin-Bauditz, S., Bingemer, H., Budke, C., Curtius, J.,

9    Danielczok, A., Diehl, K., Dreischmeier, K., Ebert, M., Frank, F.,

10   Hoffmann, N., Kandler, K., Kiselev, A., Koop, T., Leisner, T., Möhler, O.,

11   Nillius, B., Peckhaus, A., Rose, D., Weinbruch, S., Wex, H., Boose, Y.,

12   Demott, P. J., Hader, J. D., Hill, T. C. J., Kanji, Z. A., Kulkarni, G., Levin,

13   E. J. T., McCluskey, C. S., Murakami, M., Murray, B. J., Niedermeier, D.,

14   Petters, M. D., O'Sullivan, D., Saito, A., Schill, G. P., Tajiri, T., Tolbert,

15   M. A., Welti, A., Whale, T. F., Wright, T. P. and Yamashita, K.: A

16   comprehensive laboratory study on the immersion freezing behavior of

17   illite NX particles: A comparison of 17 ice nucleation measurement

18   techniques, Atmos. Chem. Phys., 15(5), doi:10.5194/acp-15-2489-2015,

19   2015a.

20   Hiranuma, N., Möhler, O., Yamashita, K., Tajiri, T., Saito, A., Kiselev, A.,

21   Hoffmann, N., Hoose, C., Jantsch, E., Koop, T. and Murakami, M.: Ice

22   nucleation by cellulose and its potential contribution to ice formation in

23   clouds, Nat. Geosci., 8(4), 273–277, doi:10.1038/ngeo2374, 2015b.

1    Lesins, G., Chylek, P., & Lohmann, U.: A study of internal and external mixing scenarios

2        and its effect on aerosol optical properties and direct radiative forcing,

3        J. Geophys. Res.-Atmos., 107, 1-12, doi:10.1029/2001jd000973, 2002.

4    Lohmann, U., and Feichter, J.: Global indirect aerosol effects: a review, Atmos. Chem.

5        Phys., 5, 715-737, doi:10.5194/acp-5-715-2005, 2005.

6    Lubin, D., and Vogelmann, A.: A climatologically significant aerosol longwave indirect

7        effect in the Arctic. Nature, 439, 453-456, doi:10.1038/nature04449, 2006.

8    Mjolsness, E.: Machine Learning for Science: State of the Art and Future Prospects,

9        Science, 293, 2051-2055, doi:10.1126/science.293.5537.2051, 2001.

10   Murphy, D. M.: The design of single particle laser mass spectrometers, Mass Spectrom.

11       Rev., 26 (2), 150–165, 2007.

12   Murphy, D. M , Middlebrook, A. M., and Warshawsky, M.: Cluster Analysis of Data

13       from the Particle Analysis by Laser Mass Spectrometry (PALMS) Instrument,

14       Aerosol Sci. Tech., 37:4, 382-391, doi:10.1080/02786820300971, 2003.

15   Niemand, M., Möhler, O., Vogel, B., Vogel, H., Hoose, C., Connolly, P., Klein, H.,

16       Bingemer, H., DeMott, P., Skrotzki, J. and Leisner, T.: A Particle-Surface-Area-

17       Based Parameterization of Immersion Freezing on Desert Dust Particles,

18       J. Atmos. Sci., 69, 3077-3092, 2012.

19   Peckhaus, A., Kiselev, A., Hiron, T., Ebert, M. and Leisner, T.: A comparative

20       study of K-rich and Na/Ca-rich feldspar ice-nucleating particles in a

21       nanoliter droplet freezing assay, Atmos. Chem. Phys., 16(18), 11477–

22       11496, doi:10.5194/acp-16-11477-2016, 2016.

1    Powers D. W.: Evaluation: From Precision, Recall and F-Factor to ROC, Informedness,

2    Markedness & Correlation, Journal of Machine Learning Technologies, 7, 1-24,

3    2007.

4    Saathoff, H., Naumann, K.-H., Schnaiter, M., Schöck, W., Möhler, O., Schurath,

5    U., Weingartner, E., Gysel, M. and Baltensperger, U.: Coating of soot and

6    (NH4)2SO4 particles by ozonolysis products of α-pinene, J. Aerosol Sci.,

7    34(10), 1297–1321, doi:10.1016/S0021-8502(03)00364-1, 2003.

8    Steinke, I., Funk, R., Busse, J., Iturri, A., Kirchen, S., Leue, M., Möhler, O.,

9    Schwartz,T., Schnaiter, M., Sierau, B., Toprak, E., Ullrich, R., Ulrich, A.,

10    Hoose, C. and Leisner, T.: Ice nucleation activity of agricultural soil dust

11    aerosols from Mongolia, Argentina, and Germany, J. Geophys. Res.

12    Atmos., doi:10.1002/2016JD025160, 2016.

13    Vogelmann, A., McFarquhar, G., Ogren, J., Turner, D., Comstock, J., Feingold, G., Long,

14    C., Jonsson, H., Bucholtz, A., Collins, D., Diskin, G., Gerber, H., Lawson, R.,

15    Woods, R., Andrews, E., Yang, H., Chiu, J., Hartsock, D., Hubbe, J., Lo,

16    C.,Marshak, A., Monroe, J., McFarlane, S., Schmid, B., Tomlinson, J. and Toto,

17    T.: Racoro Extended-Term Aircraft Observations of Boundary Layer Clouds,

18    Bull. Amer. Meteor. Soc., 93, 861-878, 2012.

19    Welti, A., Lüönd, F., Stetzer, O., and Lohmann, U.: Influence of particle size on the ice

20    nucleating ability of mineral dusts, Atmos. Chem. Phys., 9, 6929-6955,

21    doi:10.5194/acpd-9-6929-2009, 2009.

22    Zawadowicz, M. A., Froyd, K. D., Murphy, D. M. and Cziczo, D. J.: Improved

23    identification of primary biological aerosol particles using single particle

1        mass spectrometry, Atmos. Chem. Phys., doi: 10.5194/acp-2016-1119,

2        2016.

3

Atmospheric
Measurement
Techniques
Discussions

1          **Table Captions**

| Aerosol type | FIN Label | Description and/or supplier | Generation method | Sample origin | Reference |
|---|---|---|---|---|---|
| Argentinian | SDAr01 | Soil dust collected in La Pampa province, Argentina | Dry-dispersed | KIT | (Steinke et al., 2016) |
| Chinese | SDMo01 | Soil collected from Xilingele steppe, China/Inner Mongolia | Dry-dispersed | KIT | (Steinke et al., 2016) |
| Ethiopian | VSE01 | Soil collected in Lake Shala National Park, Ethiopia (collection coordinates: 7.5 N, 38.7 E) | Dry-dispersed | KIT | N/A |
| German | SDGe01 | Arable soil collected near Karlsruhe, Germany | Dry-dispersed | KIT | (Steinke et al., 2016) |
| Moroccan | DDM01 | Soil collected in a rock desert in Morocco (collection coordinates: 33.2 N, 2.0 W) | Dry-dispersed | KIT | N/A |
| Paulinenaue | N/A | Arable soil collected in Northern Germany (Brandenburg) | Dry-dispersed | KIT | N/A |
| ATD | N/A | Arizona Test Dust, Powder Technology, Inc. (Arden Hills, MN) | Dry-dispersed | MIT | N/A |
| Illite | IS03 | Illite NX (Arginotec, Germany) | Dry-dispersed | KIT | (Hiranuma et al., 2015a) |
| Fly ash | N/A | Four samples of fly ash from U.S. power plants: J. Robert Welsh Power Plant (Mount Pleasant, TX), Joppa Power Station (Joppa, IL), Clifty Creek Power Plant (Madison, IN) and Miami Fort Generating Station (Miami Fort, OH) (Fly Ash Direct, Cincinnati, OH) | Dry-dispersed | MIT | (Garimella, 2016; Zawadowicz et al., 2016) |
| Na-Feldspar | FS05 | Sodium and calcium-rich feldspar, samples provided by Institute of Applied Geosciences, Technical University of Darmstadt (Germany) and University of Leeds (UK) | Dry-dispersed | KIT | (Peckhaus et al., 2016) |
| K-Feldspar | FS01 | Potassium-rich feldspar, samples provided by Institute of Applied Geosciences, Technical University of Darmstadt (Germany) and University of Leeds (UK) | Dry-dispersed | KIT | (Peckhaus et al., 2016) |
| Agar | N/A | Agar growth medium for bacteria, Pseudomonas Agar Base (CM0559, Oxoid Microbiology Products, Hampshire, UK) | Wet-generated | KIT | N/A |
| Bacteria | PS32B74 + PFCGina01 | Two different cultures of *Pseudomonas syringae*. | Cultures grown on the agar growth medium and wet-generated | KIT | (Zawadowicz et al., 2016) |

2

3

| Cellulose | MCC01, FC01 | Microcrystalline and fibrous cellulose (Sigma Aldrich, St. Louis, MO) | Wet-generated | KIT | (Hiranuma et al., 2015b) |
|---|---|---|---|---|---|
| Hazelnut | PWW-hazelnut | Natural hazelnut pollen (GREER, Lenoir, NC) wash water | Wet-generated | KIT | (Zawadowicz et al., 2016) |
| Snomax | Snomax | Snomax, (Snomax International, Denver, CO) irradiated, desiccated and ground *Pseudomonas syringae* | Wet-generated | KIT | (Zawadowicz et al., 2016) |
| PSL | N/A | Polystyrene latex spheres (Polysciences, Inc. Warrington, PA), various sizes | Wet-generated | MIT | N/A |
| Soot | CAST minOC or maxOC | CAST soot | miniCAST flame soot generator (manufactured by Jing Ltd Zollikofen, Switzerland) | KIT | (Henning et al., 2012) |
| SOA | SOA | Secondary organic aerosol | Ozonolysis of $\alpha$-pinene | KIT | (Saathoff et al., 2003) |
| K-Feldspar cSA | FS01cSA or FS04cSA | Potassium-rich feldspar (as above) coated with sulfuric acid (SA). | Sulfuric acid incrementally added to the chamber filled with K-feldspar to achieve coatings | KIT | (Saathoff et al., 2003) |
| K-Feldspar cSOA | FS04cSOA | Potassium-rich feldspar (as above) coated with secondary organic aerosol (SOA, as above). | Sulfuric acid incrementally added to the chamber filled with K-feldspar to achieve coatings | KIT | (Saathoff et al., 2003) |

1

2    Table 1. Description of aerosol types used in training data set.

3

4

| Aerosol Type | | | | Broad Categories | | | |
|---|---|---|---|---|---|---|---|
| Negative | | Positive | | Negative | | Positive | |
| ion | feature | ion | feature | ion | feature | ion | feature |
| 35 | $^{35}Cl^-$ | 23 | $Na^+$ | 35 | $^{35}Cl^-$ | 23 | $Na^+$ |
| 25 | $C_2H^-$ | 59 | $Co^{+(1)}/CaF^+/C_2H_2OOH^+$ | 26 | $CN^-/C_2H_2^-$ | 59 | $Co^{+(1)}/CaF^+/C_2H_2OOH^+$ |
| 24 | $C_2^-$ | 39 | $^{39}K^+$ | 46 | $NO_2^-$ | 44 | $SiO^+/COO^+/^{44}Ca^+/AlOH^+$ |
| 57 | $C_2OOH^-$ | 12 | $C^+$ | 1 | $H^-$ | 39 | $^{39}K^+$ |
| 59 | $C_2H_2OOH^-/AlO_2^-$ | 24 | $C_2^+$ | 57 | $C_2OOH^-$ | 28 | $Si^+/CO^+$ |
| 43 | $HCN^-/AlO^-$ | 41 | $^{41}K^+/C_3H_5^+$ | 59 | $C_2H_2OOH^-/AlO_2^-$ | 41 | $^{41}K^+/C_3H_5^+$ |
| 1 | $H^-$ | 204-208 | Pb region ($^{204}Pb$, $^{206}Pb$, $^{207}Pb$ and $^{208}Pb$) | 45 | $COOH^-$ | 54 | $^{54}Fe^+$ |
| 26 | $CN^-/C_2H_2^-$ | 27 | $Al^+/C_2H_3^+$ | 42 | $CNO^-/C_2H_2O^-$ | 56 | $Fe^+/CaO^+$ |
| 46 | $NO_2^-$ | 44 | $SiO^+/COO^+/^{44}Ca^+/AlOH^+$ | 43 | $HCN^-/AlO^-$ | 27 | $Al^+/C_2H_3^+$ |
| 16 | $O^-$ | 57 | $^{57}Fe^+/CaOH^+/C_3H_4OH^+$ | 16 | $O^-$ | 45 | $SiOH^+/COOH^+$ |
| 17 | $OH^-$ | N/A | aerodynamic diameter | 73 | $C_2O_3H^-/C_3H_2OOH_3^-$ | 66 | $Zn^+$ |
| 61 | $SiO_2H^-/^{29}SiO_2^-/C_5H^-/CHO_3^-$ | 83 | $H_3SO_3^-/C_4H_2OOH^+$ | 63 | $PO_2^-$ | 57 | $^{57}Fe^+/CaOH^+/C_3H_4OH^+$ |
| 63 | $PO_2^-$ | 87 | $^{87}Rb^+/CaPO^+$ | 60 | $SiO_2^-/C_5^-/CO_3^-/AlO_2H^-$ | 87 | $^{87}Rb^+/CaPO^+$ |
| 19 | $F^-/H_3O^-$ | 13 | $CH^+$ | 15 | $NH^-/CH_3^-$ | 85 | $^{85}Rb^+$ |
| 76 | $SiO_3^-$ | 66 | $Zn^+$ | 24 | $C_2^-$ | 83 | $H_3SO_3^+/C_4H_2OOH^+$ |
| 77 | $SiO_3H^-/^{29}SiO_3^-$ | 28 | $Si^+/CO^+$ | 76 | $SiO_3^-$ | 24 | $C_2^+$ |
| 79 | $PO_3^-$ | 85 | $^{85}Rb^+$ | 32 | $O_2^-$ | 204-208 | Pb region ($^{204}Pb$, $^{206}Pb$, $^{207}Pb$ and $^{208}Pb$) |
| 60 | $SiO_2^-/C_5^-/CO_3^-/AlO_2^-$ | 72 | $FeO^+/CaO_2^+$ | N/A | aerodynamic diameter | 40 | $Ca^+$ |
| 45 | $COOH^-$ | 54 | $^{54}Fe^+$ | 71 | $C_3H_2OOH^-$ | 153 | $^{137}BaO^+$ |
| N/A | aerodynamic diameter | 82 | $ZnO^+$ | 50 | $C_4H_2^-$ | N/A | aerodynamic diameter |

(1) Contamination

Table 2. Features rankings for differentiation of particles between labels and between broad categories in positive and negative ion modes. See text for additional details.
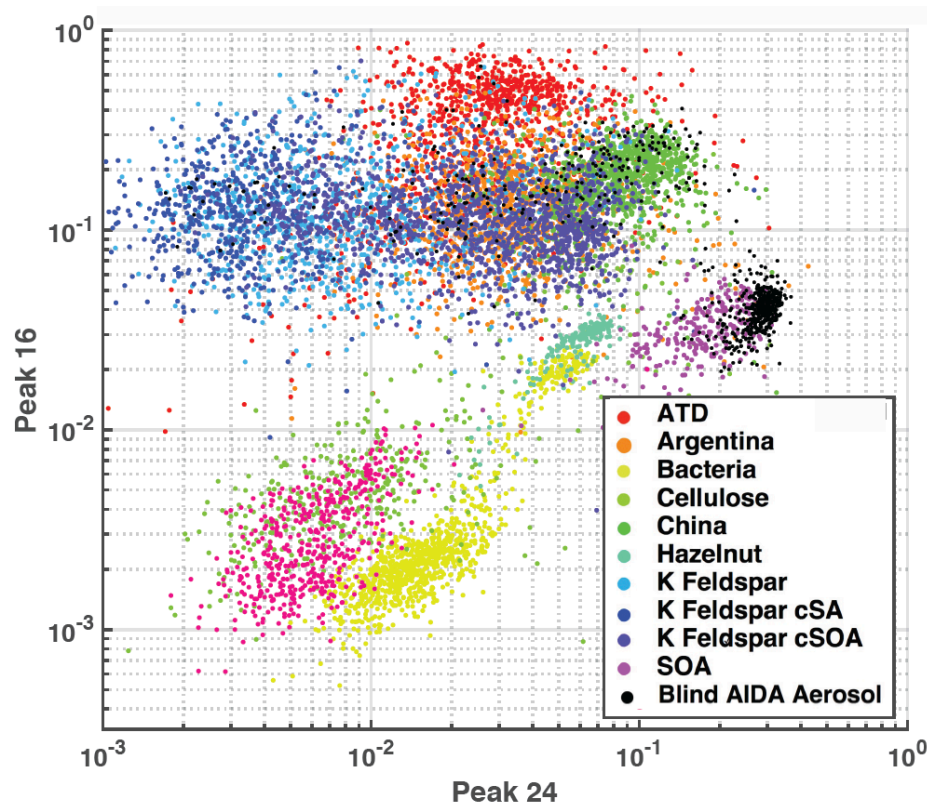
Atmospheric
Measurement
Techniques
Discussions

Open Access

28

| Category | Negative | Postive |
|---|---|---|
| Fertile Soil | 0.89 | 0.84 |
| Mineral/Metallic | 0.93 | 0.97 |
| Biological | 1.00 | 1.00 |
| Other | 0.93 | 0.96 |

| Category | Negative | Postive |
|---|---|---|
| Fertile Soil | $0.022 \pm 0.021$ | $0.032 \pm 0.031$ |
| Mineral/Metallic | $0.017 \pm 0.031$ | $0.006 \pm 0.013$ |
| Biological | 0.000 | $0.001 \pm 0.003$ |
| Other | $0.025 \pm 0.075$ | $0.010 \pm 0.029$ |

Table 3. Model performance by category and ion mode on a population consisting entirely of aerosols within that category. Left: Average classification accuracy where 1.0 = 100% precision (Powers, 2007). Right: mean and standard deviations of misclassification.

Atmospheric
Measurement
Techniques
Discussions

1    **Figure**



2

3    Figure 1: Aerosol training data plotted as feature area 16 ($O^-$) verses area 24 ($C_2^-$). Axes

4    represent peak areas normalized to total signal obtained from PALMS (i.e., 1 = 100% of

5    signal). This illustrates simple 2-dimensional clustering of aerosols from the training data

6    set by type. Co-plotted are ~500 randomly drawn spectra from the AIDA blind

7    experiment, which were known to be a subset of the training data aerosols.
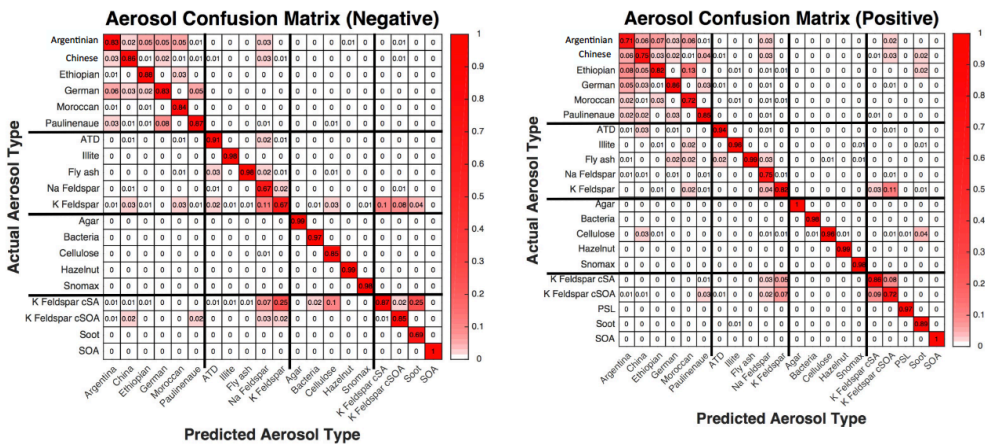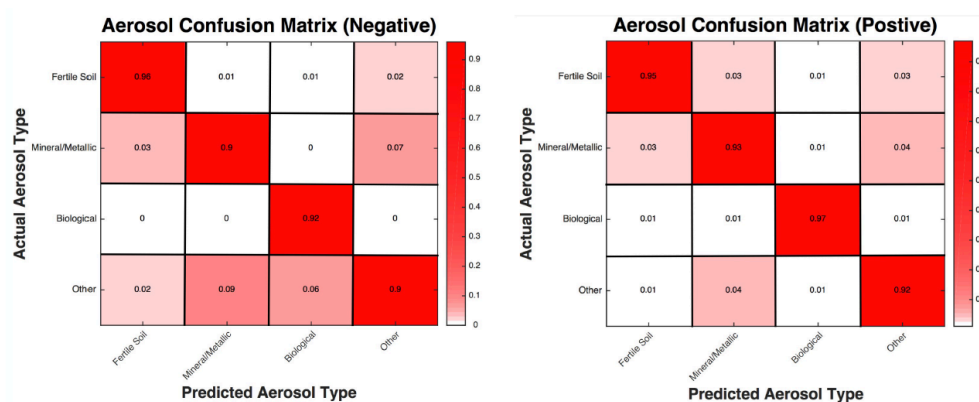
8

9

Atmospheric
Measurement
Techniques
Discussions

Open Access

1



2    Figure 2. Schematic of decision tree classification for a single aerosol spectrum. From

3    left to right, a mass spectrum is normalized with respect to total ion current, forming the

4    elements of normalized feature vector X. A trained decision tree then applies a series of

5    tests to a discreet number of peaks in order to arrive at a categorical aerosol prediction

6    (the leaves).

7

1

Figure 3. Column-normalized confusion matrices showing fraction of aerosols labeled as j that belong to i, where i and j are row and column indices, respectively. Confusion matrices are determined from training data of known origin and are used to compute probability distributions. Aerosol types (Table 1.) are grouped into four broad categories delineated by the bold horizontal and vertical bars. From top to bottom or left to right: fertile soils, mineral/metallic, biological, and other. Classification accuracy, the average probability of a correct aerosol prediction across all labels, is computed by averaging diagonal matrix elements. For all aerosol types, the accuracy is 88% in positive ion mode and 86% in negative ion mode.

Atmospheric
Measurement
Techniques
Discussions

Open Access

1

2    Figure 4. Column-normalized confusion matrices for the broad categorization of aerosols

3    following the convention in Figure 3. For all aerosol categories, the accuracy is 94% in

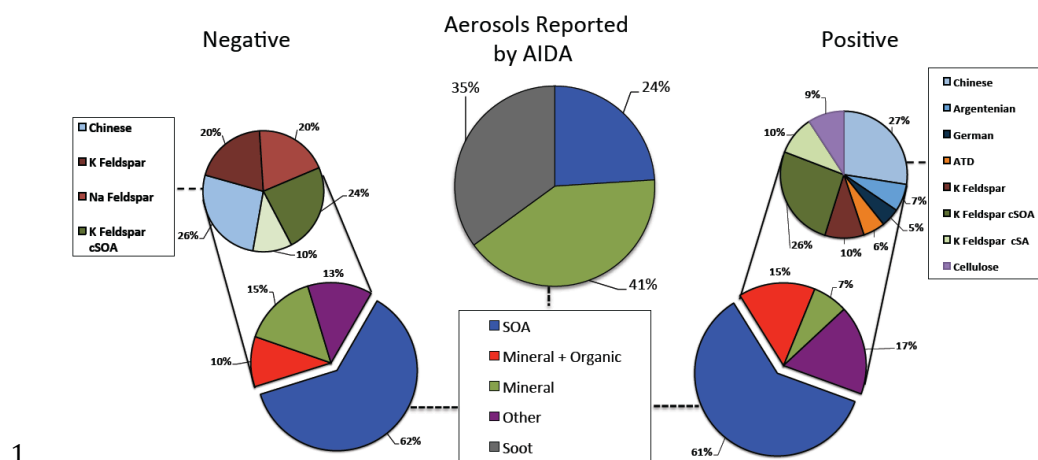4    positive ion mode and  92% in negative ion mode.

5

6

1

Figure 5. Model predictions of ~5000 aerosols sampled from the AIDA FIN01 blind mixture which was known to be a subset of the training data. Top middle : aerosol types input to the chamber for the blind mixture. Model predictions are shown for negative and positive ion mode on the left and right, respectively. Bottom: broad categories. Top: breakout by aerosol type of the non-SOA categories above the 1% level. Notes (1) the soot in the blind mixture was known to be below the instrument detection limit and therefore is not expected to be found in the data, (2) coagulation of SOA and mineral dust, which occurred after aerosol input to the chamber, appears as the mineral + organic category.