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] now at ACME AtronOmatic, LLC, Portland, OR, United States

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

11 National 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

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*. Our primary focus surrounds the growing of random forests

9   using feature selection to reduce dimensionality, and the evaluation of trained models

10  with confusion matrices.  In addition to classifying ~20 unique, but chemically-similar,

11  aerosol types, models were also created to differentiate aerosol within four broader

12  categories: fertile soils, mineral/metallic particles, biological, and all other aerosols.

13  Differentiation was accomplished using ~40 positive and negative spectral features. For

14  the broad categorization, machine learning resulted in a classification accuracy of ~93%.

15  Classification of aerosols by specific type resulted in a classification accuracy of ~87%.

16  The 'trained' model was then applied to a 'blind' mixture of aerosols which was known

17  to be a subset of the training set. Model agreement was found on the presence of

18  secondary organic aerosol, coated and uncoated mineral dust and fertile soil.

19

20  **1. Introduction**

1       Following the introduction of random forests in the 1990s, recent developments in

2 deep learning and neural networks have triggered a renewed interest in machine learning.

3 This has led to the development of numerous easy-to-use, freely-available, open-source

4 packages in popular programming languages like Python, and these tools are becoming

5 increasing used in academia and industry. While random forests have been used for

6 complex classification and regression analysis in various fields, studies that employ

7 random forests in aerosol mass spectrometry remain sparse. Utilizing these tools, the

8 primary purpose of our study is to introduce a framework for growing random forests,

9 reducing dimensionality, ranking chemical features, and evaluating performance using

10 confusion matrices. Such properties are desirable for SPMS studies, where input

11 variables can become redundant and interpretability is more limited with methods such as

12 cluster analysis and neural networks. Analysis techniques such as those falling out of

13 recent artificial intelligence research can prove useful for helping to tease out the subtle

14 yet significant impact that aerosol chemistry has on the climate system.

15       Atmospheric aerosols impact clouds and the Earth's radiative budget. A lack of

16 understanding of aerosol composition therefore contributes to uncertainty in

17 determination of both anthropogenic and natural climate forcing [Boucher et al., 2013;

18 Lohmann and Feichter, 2005]. Aerosols directly affect atmospheric radiation by

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

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

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

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

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

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

2    morphology [Lohmann and Feichter, 2008]. Characterization of aerosol properties plays a

3    vital role in understanding weather and climate.

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

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

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

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

8    aerosol mixtures and characterizing both semi-volatile (e.g. organics and sulfates) and

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

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

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

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

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

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

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

16         Chemical composition of an individual atmospheric aerosol particle is a complex

17   interplay between its primary composition at the source (i.e. dust, biogenic organic,

18   anthropogenic organic, soot, etc.) and its atmospheric processing up to the time of

19   detection. Atmospheric processing can include a combination of coating with secondary

20   material, coagulation and cloud processing. Even different primary aerosol types can

21   have similar mass spectral markers. For example, fly ash, mineral dust and bioaerosol can

22   all contain strong phosphate signal [Zawadowicz et al., 2017]. Secondary material is

23   often difficult to differentiate from primary material, but even minor compositional

1    changes can be atmospherically important. As one example, mineral dusts are known to

2    be effective at nucleating ice clouds; however, despite minor addition of mass,

3    atmospherically processed mineral dust is less suitable for ice formation [Cziczo et al.,

4    2013]. As a second example, ice nucleation in mixed-phase clouds has been suggested to

5    be predominantly influenced by feldspar, a single component among the diverse

6    mineralogy of atmospheric dust [Atkinson et al., 2013]. Using current SPMS data

7    analysis approaches, it is difficult to detect these minor yet important compositional

8    differences and new robust and generalizable analysis techniques are critical.

9        We show that supervised training with random forests can differentiate aerosols in

10    SPMS data more accurately than simpler approaches. Various clustering methods have

11    been used to group aerosol types [Murphy et al., 2003; Gross et al., 2008] but these

12    algorithms are known to struggle with chemically-similar aerosols as they do not

13    incorporate known particle labels in the training process. Such 'unsupervised' clustering

14    algorithms automatically group unlabeled data points in feature space, in this case mass

15    spectral signals. For the purposes of setting broad aerosol categories, which are

16    chemically similar and easily separable in feature space, clustering is the simpler tool and

17    the data easier to interpret. For identifying new or potentially unexpected atmospheric

18    aerosols, such properties are desirable; however, the advantages of clustering greatly

19    diminish when considering similar particle types that overlap in feature space. A

20    limitation often encountered is the need to manually reduce the number of final clusters

21    due to grouping of mathematically-similar yet chemically-distinct aerosols [Murphy et al

22    2003]. Fertile soils, for instance, are often grouped into a single category despite

23    different sources and atmospheric histories. Clustering algorithms should therefore be

considered as a tool to use alongside supervised classification. The latter may be used to further explore unique aerosol types or verify manually labeled clusters with higher precision. Furthermore, the ensemble approach presented here also produces interpretable variable rankings and probabilistic predictions that assist in addressing measurement uncertainty. The choice of supervised or unsupervised machine learning will depend on the researcher's use-case, and each method has unique advantages and disadvantages. We note a limitation of the random forest approach - and for supervised learning in general - is the inability to classify aerosol types outside of the training set. The ability of a random forest to characterize ambient atmospheric datasets, therefore, will strongly depend on which aerosols are contained within the training set.

In this study, we demonstrate the capabilities of random forests to automatically differentiate particles on the basis of chemistry and size. The resulting model can capture minor compositional differences between aerosol mass spectra. By testing predictions using an independent, or 'blind', dataset, we illustrate the feasibility of combining on-line analysis techniques such as SPMS with machine learning to infer the behavior and origin of aerosols in the laboratory and atmosphere.

## 2. Methodologies

### 2.1 PALMS

The Particle Analysis by Laser Mass Spectrometry (PALMS) instrument was employed for these studies. PALMS has been described in detail previously [Cziczo et al. 2006]. Briefly, the instrument samples aerosol particles in the size range from ~200 to ~3000 nm using an aerodynamic lens inlet into a differentially-pumped vacuum region.

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

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

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

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

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

6    refractory and semi-volatile components and allows analysis of all chemical components

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

8

9    **2.2 Dataset**

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

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

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

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

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

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

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

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

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

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

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

21   ~50,000 spectra are considered with each spectrum containing 512 possible mass peaks

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

23   sampling period, where AIDA was filled with an unknown number of aerosol types

1  known to be from the training set (i.e., for which spectra had already been acquired) but

2  (*a priori*) of unknown size, specific types and at unknown concentrations.

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

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

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

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

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

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

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

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

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

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

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

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

15  **2.3 Aerosol Classification**

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

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

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

19  classification methods, machine learning determines parameters that partition the data set.

20  To form X, mass spectra are converted to dimensional vectors normalized to the total ion

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

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

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

1    variables for the machine learning model.

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

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

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

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

6    type. A random forest was used to generate predictions of aerosol type. A single decision

7    tree is a statistical decision model that performs classification based on a series of

8    comparisons relating a variable $X_i$ (in this case a normalized mass peak in X) to a learned

9    threshold value [Breiman, 2001]. Represented as an algorithmic tree, a binary decision

10   tree consists of a hierarchy of nodes where each node connects via branches to two other

11   nodes deeper in the tree. At each node, one of the two branches is taken based on whether

12   a normalized peak $X_i$ is greater or less than a threshold value. Each branch leads to

13   another node where a different test is performed. After a series of tests, one at each node,

14   a class is assigned to a given sample; these are the so-called 'leaves'. Figure 2 illustrates

15   the classification model for a single decision tree.

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

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

18   10 to 3000, the set of reachable leaves and possible labels is 1 and the decision tree

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

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

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

22   each polarity, were generated to classify aerosols. These are hereafter referred to as the

23   'positive' and 'negative classification algorithms'.

1    **2.4 Random Forests**

2        A random forest is an ensemble of decision tree classifiers where each classifer

3    independently labels an unknown spectrum vector X. To make a final prediction of

4    aerosol type, trees within an ensemble 'vote' on a classification label. Each vote has

5    equal weight and the spectrum is assigned to the majority choice. Each tree within an

6    ensemble is independently grown on a subset of the training data so that a commonly

7    voted label implies a higher certainty. Adding members to an ensemble increases the

8    robustness of a classification model by providing alternative hypotheses and is therefore

9    preferable to single classifiers.

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

11    independently grown during training.  A total of $k$ trees, with $k = 110$, were grown using a

12    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 original data, leaving a

15    portion of the training set unsampled. The unsampled data for each tree, known as 'out-

16    of-bag' observations, are recorded and later provide a means to assess classification error

17    for the forest. To determine model error, predictions are made for each point in the

18    dataset using only the subset of trees that did not use the point for training. Each training

19    point is left out at least once. This is analogous to making predictions with a separately

20    trained forest that did not observe the point and prevents testing with the same data used

21    for training.

22        Given a bootstrap sample, a tree is grown by sequentially creating tests that

1 maximize the separation between classes in parameter space. A test is created by

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

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

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

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

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

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

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

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

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

11     The number of variables per split is chosen to be 11 and the number of trees is

12 110. The optimal model was determined by enumerating combinations of these

13 parameters on a coarse grid and selecting the values that produce the lowest test error.

14 Model behavior is primarily sensitive to the number of variables per split, and shows

15 weak dependence on the number of trees and number of input variables beyond small

16 values. As the number of variable splits increases, error decreases exponentially to a local

17 minimum before again rising due to over fitting. Alternatively, as the number of trees is

18 increased the error asymptotes to some nonzero value, a known characteristic of random

19 forests where test error converges to the generalization error. The models were trained

20 with the Python 2.7 Scikit-learn module on a MacBook Pro with 16 GB 1600 MHz

21 DDR3 memory and a 2.5 GHz Intel Core i7 processor. A typical random forest model

22 took about 5-10 seconds to train, and we found a linear relationship between runtime and

23 both the number of trees and variables per split.

1       Overall, the generalizability and robust performance of random forests is owed

2    significantly to the series of random statistical procedures used to construct such models.

3    An ensemble classifier reduces variability by averaging predictions over a series of

4    independently trained models, and bagging introduces additional randomness by

5    producing "perturbed" versions of the original data via random sampling of input data.

6    The randomness used in constructing forests, both in bagging the training set and

7    choosing variable splits, work to decorrelate the output of each tree even as the inputs

8    become correlated [Breiman, 2001]. As the number of trees increases, the law of large

9    numbers guarantees a convergence of the out-of-bag error to the generalization error.

10    **2.5 Dimensionality Reduction and Chemical Feature Selection**

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

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

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

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

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

16    differentiation. Feature ranking was algorithmically determined by comparing the

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

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

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

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

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

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

23    for correct classification) for the entire set of aerosol types. The columns at right list

rankings when aerosol types are grouped into the broad, chemically similar, categories. A final ranking was determined by sequentially adding variables and observing classification performance response. All variables preceding two e-foldings in classification error were maintained in the final model. Both the specific aerosol type and broad aerosol category models were retrained using this subset of the initial variables, listed in Table 2.

**2.5 Comparison to Euclidean Distance Classifier**

To access relative model performance, we contrast the results with a simple classifier that compares unseen aerosols to a set of class mean vectors. Using the Euclidean distance metric, the unknown aerosol is assigned to the nearest class. This simple baseline classifier helps put results in the context of machine learning techniques that rely on distance-based metrics such as k-means and hierarchical clustering. K-means clustering attempts to divide the data points into k distinct clusters, representing spectra as vectors. Using Euclidean distance, the standard algorithm assigns points to centroids, or clusters, which are essentially mean vectors representing the average of all points in the cluster. Assuming perfect convergence of k-means clustering, where k is the number of aerosol classes, each cluster represents the mean of aerosol in that class. The random forest results below demonstrate many areas of improvement over the simple classifier.

**3. Results**

**3.1 Confusion Matrices and Probabilistic Model Performance**

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

2    the model prediction with the true aerosol type or broad category [Powers, 2007], and can

3    be understood as a contingency table matching model predictions to true labels.

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

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

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

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

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

9    basis of classification is separation of physical quantities, misclassifications result from

10   similarity in mass peaks and their ion abundance between aerosol types. This is most

11   easily visualized as overlapping clusters in the simple two dimensional space in Figure 1.

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

13   the confusion matrix, which represent the fraction of aerosol in column j labeled

14   correctly. The classification accuracy (*a*) is given by averaging diagonal elements of P. A

15   perfect classification model produces the identity matrix, as all data points are classified

16   correctly 100% of the time. For example, in the positive confusion matrix, SOA and Agar

17   growth medium are correctly labeled in the test set 100% of the time. Barring element

18   truncation, all columns 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. The precision score [Powers, 2007] captures the classification behavior for

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

14 within that category:

15 $$\text{Precision Score(L)} = \frac{1}{|L|} \sum_{i=j}^{|L|} P(i \in L, j \in L) \qquad (3)$$

16

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

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

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

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

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

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

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

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

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

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

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

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

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

6    labeling rate.

7        Across all categories, the random forest shows improvements over the Euclidean

8    classifier in terms of both accuracy and precision. Figure 4 directly compares confusion

9    matrices for the two methods, revealing overall accuracy improvements of at least 20%.

10   The largest improvements are in the fertile soil and other category, where accuracy rises

11   between 20% and 39% with the random forest. Computing the full confusion matrix for

12   the Euclidean technique (as in figure 3) reveals similar results, with far more frequent

13   mislabeling between fertile soils as well as coated/uncoated particles than our approach.

14   These results reinforce the fact that chemically-similar aerosols which overlap in feature

15   space will often be grouped together when using a single, distance-based classifier. The

16   improvement from random forests is likely a result of a) the ensemble approach, which is

17   known to produce better generalizability than single classifiers and b) the tendency of

18   aerosols   with   similar   chemical   properties   and   atmospheric   effect   to   appear

19   mathematically distinct with a distance metric.

20        Beyond classification, the obtained variable rankings alone provide interesting

21   insights into the dataset. It is noteworthy that while most of the features are logical

22   differentiators of the aerosol types investigated in FIN01 there were also surprises. One

23   example is $59^+$ (cobalt), determined to be one of the most important features for

1    differentiation. Further investigation determined this material was associated with

2    tungsten carbide contaminant from dry powder dispersion equipment used on some

3    samples. The contamination affected feldspar samples used during the second half of the

4    AIDA measurements in particular. This serves to illustrate the lack of *a priori* judgment

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

6    contamination identification).

7

8    **3.2 Characterization of Blind Data**

9    As part of the FIN01 workshop, it was known that an unknown number of aerosol

10   types from Table 1 were aerosolized into the ADIA chamber at unknown size and relative

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

12   ~25,000 spectra. After data analysis, the aerosol types and relative abundances were

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

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

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

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

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

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

19   atmosphere.

20   Figure 5 illustrates the fractional percentages for each aerosol category. Because

21   SOA was nearly always labeled correctly (Figure 3), the remaining aerosols are

22   considered separately using the full set of candidate aerosol labels. Both positive and

23   negative models arrived at similar results, with inconsistencies primarily associated with

1    the presence of trace fertile soils and mineral dust / fly ash particles. The positive

2    algorithm identifies ~2-4% of the AIDA population as each Argentinean soil, German

3    soil, ATD, and cellulose whereas the frequency of these aerosols was too low to consider

4    in the negative. Alternatively, the negative model estimates Na-Feldspar at ~14% of the

5    total population, a label not identified by the positive algorithm. This discrepancy can be

6    explained by the 1% selection criterion for aerosols present in the population. Fertile soils,

7    ATD, and cellulose frequently accumulate error along rows in the full positive confusion

8    matrix, indicating frequent confusion with other categories (Figure 3). Furthermore, with

9    the observed misclassification rates ranging ~1-4%, it is expected that these aerosol

10   labels are false positives. The negative model offers an alternative hypothesis, suggesting

11   these miscellaneous aerosols are Na-feldspar. Since there is significant model agreement

12   on the percentages of SOA, K-Feldspar, and coated feldspars, this part of the blind

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

14   aerosol labels, more credence is lent to the negative classification algorithm on the basis

15   of improved precision for fertile soils.

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

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

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

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

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

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

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

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

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

3    AIDA chamber during the course of the blind experiment. This may also explain some

4    indications of fertile soils, which are known to be mixtures of mineral and organic

5    components. The training data set did not contain coagulated SOA and mineral dust but

6    did include SOA-coated K-Feldspar, which explains the identification.

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

8    results are largely consistent with the presence of coagulated organics and minerals and

9    the known uncertainties highlighted by the confusion matrices discussed previously.

10   Given the presence of any single mineral dust, some confusion with fertile soils, SA

11   coated Feldspar, and Na-Feldspar is expected (Figure 3). Moreover, as discussed

12   previously [Gallavardin et al., 2008], AIDA backgrounds are not completely particle-free.

13   During the FIN01 study, contamination particles from previous test aerosol were

14   frequently observed as background and they could also be the origin of some low-

15   concentration particles matching fertile soil chemistry.

16

## 4. Conclusions and Future Work

18   This study lays out a framework for training and implementing random forests on

19   SPMS data, with a focus on dimensionality reduction and the evaluation of model

20   performance with confusion matrices. A key benefit to the proposed method is chemical

21   feature selection, which allows researchers to identify potentially important chemical

22   markers between arbitrary groups of aerosols or identify sources of contamination.

1    Additionally, the approach allows for differentiation of aerosols within a SPMS dataset,

2    augmenting existing tools and reducing the need for a qualitative comparison between

3    mass spectra. Across a representative sample of possible aerosol types, the behavior of

4    each algorithm predictably allows users to infer the presence or absence of specific

5    aerosols and quantify aerosol abundance. Machine learning is automated and the output

6    of the model must then be informed by human knowledge of aerosol chemistry. Machine

7    learning should therefore be considered as an additional tool to interpret mass spectra to

8    better distinguish aerosols with unique properties in terms of atmospheric chemistry,

9    biogenic cycles, and population health.

10        The random forest classification framework described here may be generalized to

11    any instrument, or set of instruments, capable of collecting physical and chemical

12    information that distinguishes particles. Although the method described here is applied to

13    a stand-alone SPMS and tested with a set of 'blind' data, ancillary laboratory or field data

14    can be integrated to expand the data set. The success of these algorithms is data-

15    dependent, where better performance is expected for instruments that provide more, and

16    more quantitative, analysis of the aerosol properties. Although the algorithms

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

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

19    **Acknowledgements**

3

## 4    References

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

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

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

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

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

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

11       doi:10.1038/nature12278, 2013.

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

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

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

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

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

17       657, 2013.

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

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

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

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

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

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

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

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

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

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

7       doi:10.1126/science.1234145, 2013.

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

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

10      1772–1775, 2001.

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

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

13      56-63, doi:10.1016/j.ijms.2008.04.031, 2008.

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

15      Single Particle Laser Mass Spectrometry Applied to Differential Ice Nucleation

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

17      10.1080/02786820802339538, 2008.

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

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

20      (in prep)

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

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

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

760-769, 2008.

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

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

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

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

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


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

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

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

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

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

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

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

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

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

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

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

comprehensive laboratory study on the immersion freezing behavior of

illite NX particles: A comparison of 17 ice nucleation measurement

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

2015a.

Hiranuma, N., Möhler, O., Yamashita, K., Tajiri, T., Saito, A., Kiselev, A., Hoffmann, N., Hoose, C., Jantsch, E., Koop, T. and Murakami, M.: Ice nucleation by cellulose and its potential contribution to ice formation in clouds, Nat. Geosci., 8(4), 273–277, doi:10.1038/ngeo2374, 2015b.

Lesins, G., Chylek, P., & Lohmann, U.: A study of internal and external mixing scenarios and its effect on aerosol optical properties and direct radiative forcing, J. Geophys. Res.-Atmos., 107, 1-12, doi:10.1029/2001jd000973, 2002.

Lohmann, U., and Feichter, J.: Global indirect aerosol effects: a review, Atmos. Chem. Phys., 5, 715-737, doi:10.5194/acp-5-715-2005, 2005.

Lubin, D., and Vogelmann, A.: A climatologically significant aerosol longwave indirect effect in the Arctic. Nature, 439, 453-456, doi:10.1038/nature04449, 2006.

Mjolsness, E.: Machine Learning for Science: State of the Art and Future Prospects, Science, 293, 2051-2055, doi:10.1126/science.293.5537.2051, 2001.

Murphy, D. M.: The design of single particle laser mass spectrometers, Mass Spectrom. Rev., 26 (2), 150–165, 2007.

Murphy, D. M , Middlebrook, A. M., and Warshawsky, M.: Cluster Analysis of Data from the Particle Analysis by Laser Mass Spectrometry (PALMS) Instrument, Aerosol Sci. Tech., 37:4, 382-391, doi:10.1080/02786820300971, 2003.

Niemand, M., Möhler, O., Vogel, B., Vogel, H., Hoose, C., Connolly, P., Klein, H., Bingemer, H., DeMott, P., Skrotzki, J. and Leisner, T.: A Particle-Surface-Area-Based Parameterization of Immersion Freezing on Desert Dust Particles, J. Atmos. Sci., 69, 3077-3092, 2012.

Peckhaus, A., Kiselev, A., Hiron, T., Ebert, M. and Leisner, T.: A comparative study of K-rich and Na/Ca-rich feldspar ice-nucleating particles in a nanoliter droplet freezing assay, Atmos. Chem. Phys., 16(18), 11477–11496, doi:10.5194/acp-16-11477-2016, 2016.

Powers D. W.: Evaluation: From Precision, Recall and F-Factor to ROC, Informedness, Markedness & Correlation, Journal of Machine Learning Technologies, 7, 1-24, 2007.

Saathoff, H., Naumann, K.-H., Schnaiter, M., Schöck, W., Möhler, O., Schurath, U., Weingartner, E., Gysel, M. and Baltensperger, U.: Coating of soot and $(NH4)2SO4$ particles by ozonolysis products of $\alpha$-pinene, J. Aerosol Sci., 34(10), 1297–1321, doi:10.1016/S0021-8502(03)00364-1, 2003.

Steinke, I., Funk, R., Busse, J., Iturri, A., Kirchen, S., Leue, M., Möhler, O., Schwartz,T., Schnaiter, M., Sierau, B., Toprak, E., Ullrich, R., Ulrich, A., Hoose, C. and Leisner, T.: Ice nucleation activity of agricultural soil dust aerosols from Mongolia, Argentina, and Germany, J. Geophys. Res. Atmos., doi:10.1002/2016JD025160, 2016.

Vogelmann, A., McFarquhar, G., Ogren, J., Turner, D., Comstock, J., Feingold, G., Long, C., Jonsson, H., Bucholtz, A., Collins, D., Diskin, G., Gerber, H., Lawson, R., Woods, R., Andrews, E., Yang, H., Chiu, J., Hartsock, D., Hubbe, J., Lo, C.,Marshak, A., Monroe, J., McFarlane, S., Schmid, B., Tomlinson, J. and Toto, T.: Racoro Extended-Term Aircraft Observations of Boundary Layer Clouds, Bull. Amer. Meteor. Soc., 93, 861-878, 2012.

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

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

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

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

5         identification of primary biological aerosol particles using single particle

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

7         2016.

8

1 **Table Captions**

2

| Aerosol type | FIN Label | Description and/or supplier | Generation method | Sample provided by | Reference 3 |
|---|---|---|---|---|---|
| **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 (as above), suspended in nanopure water and wet-generated | KIT | (Zawadowicz et al., 2016) |
| **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). | Small amounts of sulfuric acid were incrementally added to the chamber filled with K-feldspar to achieve thin coatings, as judged from PALMS spectra | KIT | (Saathoff et al., 2003) |
| **K-Feldspar cSOA** | FS04cSOA | Potassium-rich feldspar (as above) coated with secondary organic aerosol (SOA, as above). | Small amounts of SOA were incrementally added to the chamber filled with K-feldspar to achieve thin coatings, as judged from PALMS spectra | KIT | (Saathoff et al., 2003) |

1  Table 1. Description of aerosol types used in training data set. Rows are grouped and

2  colored by broad aerosol categories in the following order: Fertile Soil, Mineral/Metallic,

3  Biological, and Other.

4

5

| 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 | $SiO_2^-/C_5^-/CO_3^-/ AlO_2H^-$ | 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.

| Category | Negative | Postive |
|---|---|---|
| Fertile Soil | 0.88 | 0.83 |
| Mineral/Metallic | 0.93 | 0.98 |
| Biological | 1.00 | 1.00 |
| Other | 0.96 | 0.93 |

| Category | Negative | Postive |
|---|---|---|
| Fertile Soil | $0.024 \pm 0.020$ | $0.035 \pm 0.033$ |
| Mineral/Metallic | $0.017 \pm 0.027$ | $0.006 \pm 0.008$ |
| Biological | $0.000$ | $0.001 \pm 0.002$ |
| Other | $0.021 \pm 0.015$ | $0.024 \pm 0.053$ |

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.

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

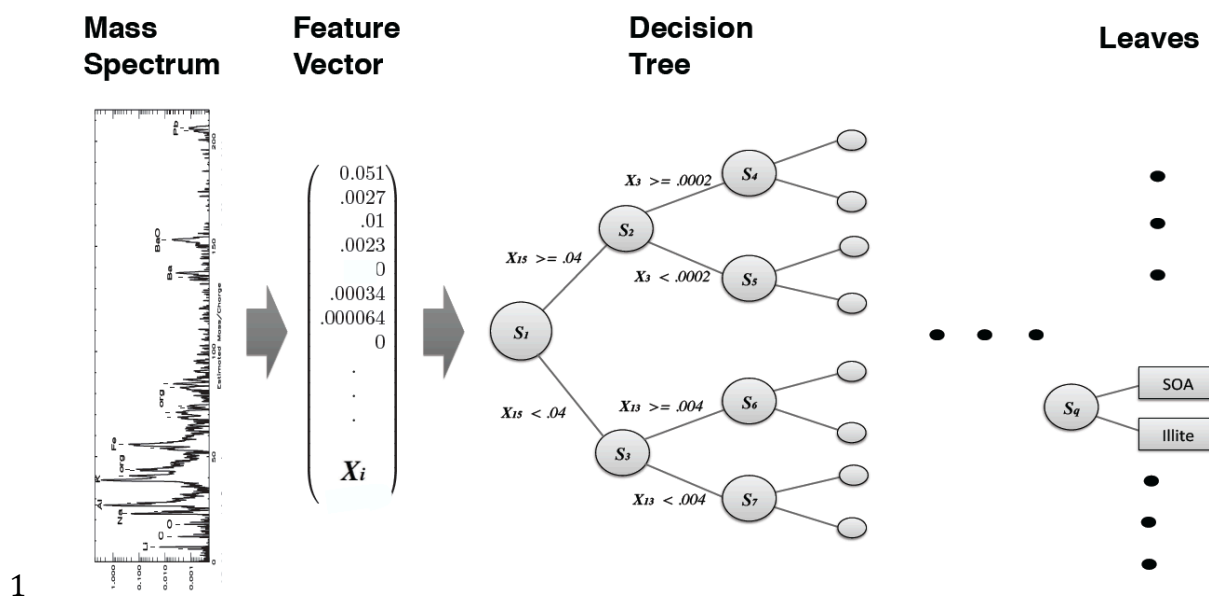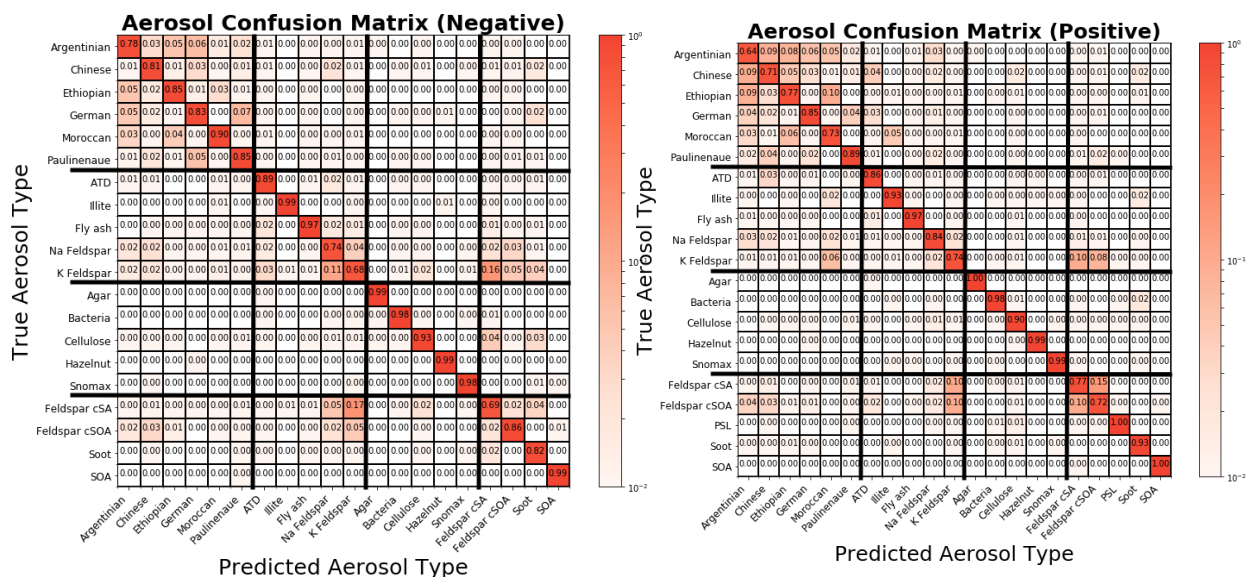7    experiment, which were known to be a subset of the training data aerosols.

8

9

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

2   Figure 3. Column-normalized confusion matrices showing fraction of aerosols labeled as

3   j that belong to i, where i and j are row and column indices, respectively. Confusion

4   matrices are determined from training data of known origin and are used to compute

5   probability distributions. Aerosol types (Table 1.) are grouped into four broad categories

6   delineated by the bold horizontal and vertical bars. From top to bottom or left to right:

7   fertile soils, mineral/metallic, biological, and other. Classification accuracy, the average

8   probability of a correct aerosol prediction across all labels, is computed by averaging

9   diagonal matrix elements. For all aerosol types, the accuracy is 87% in positive ion mode
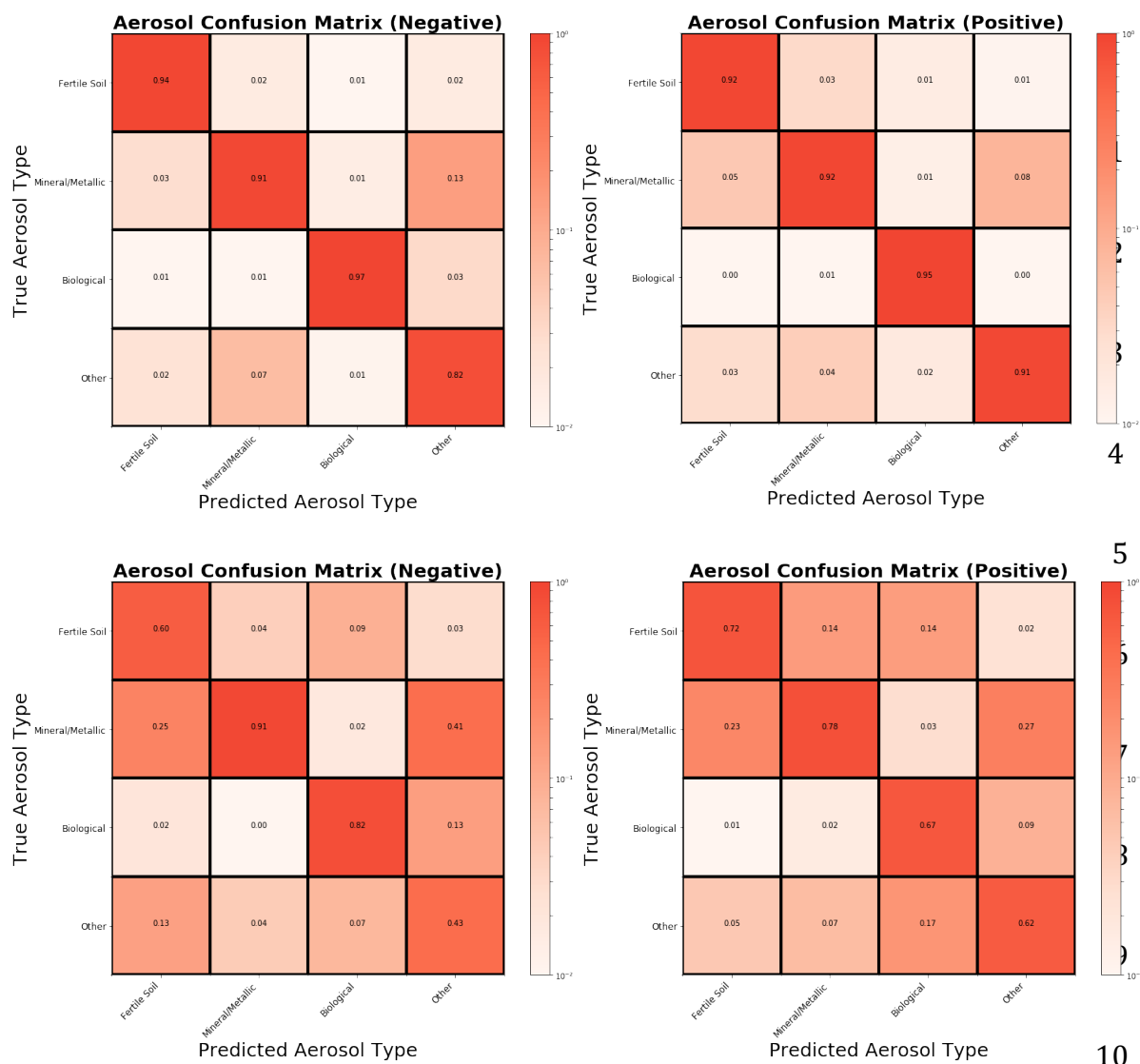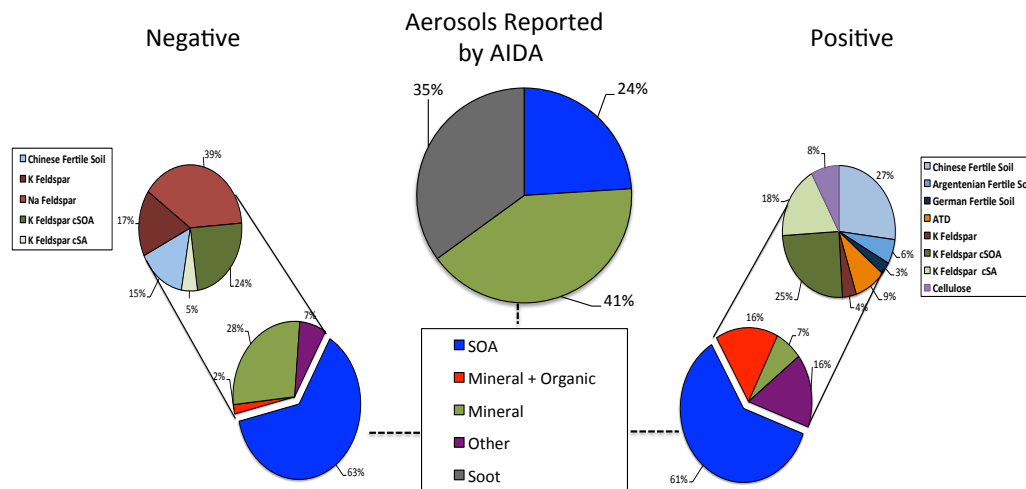
10   and  87% in negative ion mode.

11

Figure 4. Column-normalized confusion matrices for the broad categorization of aerosols following the convention in Figure 3. a.)Top row: For all aerosol categories, the random forest has an accuracy of 93% in positive ion mode and 91% in negative ion mode. b.) Bottom row: The Euclidean distance classifier has an accuracy of 70% in positive ion mode and 69% in negative ion mode

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, was often categorized as mixed mineral and organic particles or fertile soils (i.e., mixtures of mineral and organic components) considered in the training data set, (3) the aerosols types reported by AIDA do not account for PALMS transmission efficiency (see text for details).