



Supplement of

On the competition among aerosol number, size and composition in predicting CCN variability: a multi-annual field study in an urbanized desert

E. Crosbie et al.

Correspondence to: A. Sorooshian (armin@email.arizona.edu)

The copyright of individual parts of the supplement might differ from the CC-BY 3.0 licence.

Explanation of size distribution shape clustering method

Each size distribution observation comprises 112 number concentrations, n, corresponding to the discrete bins reported by the SMPS software on a logarithmic, regularly spaced grid. Each observation is considered as a vector, \mathbf{x}_i , normalized using the root sum of squared distance (Euclidean norm) to isolate the effect of size distribution shape from the (strong) variability in total CN:

$$\mathbf{x}_{i}^{*} = \{n_{1}, n_{2}, \dots, n_{112}\}$$

 $\mathbf{x}_{i} = \mathbf{x}_{i}^{*} / \|\mathbf{x}_{i}^{*}\|_{2}$

The K-means algorithm seeks the optimal position of a predetermined number, k, of centroid vectors, \mathbf{c}_j which minimizes the residual distance of observation vectors assigned to that cluster set, C_j , from the centroid. Centroids are initially randomly assigned to observations and then temporary cluster assignments are made to the closest centroid for each observation, \mathbf{x}_i . The centroid is then recalculated based on the cluster members and the process is iterated until convergence. The method is known to sometimes converge on local optima rather than global optima in some cases and so the operation is repeated many times (>100) with different random starting observations to check for global convergence.

$$C_{j}^{(m)} = \left\{ \boldsymbol{x}_{i} : \left\| \boldsymbol{x}_{i} - \boldsymbol{c}_{j}^{(m)} \right\|_{2}^{2} \leq \left\| \boldsymbol{x}_{i} - \boldsymbol{c}_{p}^{(m)} \right\|_{2}^{2} \quad \forall p \right\}$$
$$\boldsymbol{c}_{i}^{(m+1)} = \frac{1}{|C_{j}^{(m)}|} \sum_{\boldsymbol{x}_{i} \in C_{j}^{(m)}} \boldsymbol{x}_{i}$$

The resulting cluster centroids after convergence represent the mean size distribution shape of the subset of observations attributed to that cluster. Observations show a quasi-continuous nearest neighbor density in the regions between the centroids, hence using the "hard" cluster associations (as per the traditional K-means algorithm) creates an abrupt transition for observations which lie close to the cluster boundary. Instead the cluster association is made "fuzzy" using cluster assignment weights, w_j. Weights are based on inverse squared Euclidean distance to each centroid, c_j, except where the square distance is greater than the sum of squared distances via any other centroid, c_p, in which case the weight, w_j, is zero. Weights are normalized such that w_j =1.

$$d_{i,j} = \left\| \boldsymbol{x}_{i} - \boldsymbol{c}_{j} \right\|_{2}$$
$$w_{i,j}^{*} = \begin{cases} d_{i,j}^{-2}; & d_{i,j}^{2} \leq d_{i,p}^{2} + \left\| \boldsymbol{c}_{p} - \boldsymbol{c}_{j} \right\|_{2}^{2} \\ 0; & d_{i,j}^{2} > d_{i,p}^{2} + \left\| \boldsymbol{c}_{p} - \boldsymbol{c}_{j} \right\|_{2}^{2} \end{cases}$$

$$w_{i,j} = w_{i,j}^* / \sum_j w_{i,j}^*$$

The chosen number of cluster centroids (4) was selected subjectively to balance the improvement in model performance by increased numbers of clusters with the simplicity of retaining as few as possible. The quality of the K-means model as the number of centroids is changed can be visualized by the sum of residual distances between observations and closest cluster centroid. This is shown in Figure S1. The cluster centroids were also identified using only the data from the summer seasons (PM and M) since these seasons contain the largest range of variability in the size distribution shape. Repeating the K-means algorithm using the entire dataset finds the clusters common in winter (FF and WN) but finds a hybrid of the more extreme summertime clusters (N and CC) with the FF and WN clusters, respectively. Since the goal is to be able to isolate features, the clusters based on the summer "training" data are retained. Associations based on the fuzzy logic described above are then assigned to all data.



Figure S1: Performance of the K-means method as a function of the number of centroids, k. The residual is measured as the Euclidean norm of the vector difference between the observation and its allocated centroid and the mean is generated over all "training" observations. For comparison, the performance of the K-means model is shown for uncorrelated random noise and highlights the latent structure of the observations.