# Data-driven clustering of rain events: microphysics information derived from macro scale observations

Mohamed Djallel Dilmi[1], Cécile Mallet[1], Laurent Barthes[1], Aymeric Chazottes[1]

[1]LATMOS-CNRS/ UVSQ/ UPSay, 11 boulevard d'Alembert, 78280 Guyancourt, France

Corresponding author: Laurent Barthes (laurent.barthes@latmos.ipsl.fr)

**Abstract**. Rain time series records are generally studied using rainfall rate or accumulation parameters, which are estimated for a fixed duration (typically 1 min, 1 hour or 1 day). In this study we use the concept of "rain events". The aim of the first part of this paper is to establish a parsimonious characterisation of rain events, using a minimal set of variables selected among those normally used for the characterization of these events. A methodology is proposed, based on the combined use of a Genetic Algorithm (GA) and Self Organising Maps (SOM). It can be advantageous to use a SOM, since it allows a high dimensional data space to be mapped onto a two-dimensional space while preserving, in an unsupervised manner, most of the information contained in the initial space topology. The 2D maps obtained in this way allow the relationships between variables to be determined and redundant variables to be removed, thus leading to a minimal subset of variables. We verify that such 2D maps make it possible to determine the characteristics of all events, on the basis of only five features (the event duration, the peak rain rate, the rain event depth, the standard deviation of the rain rate event, and the absolute rain rate variation of the order of 0.5). From this minimal subset of variables, hierarchical cluster analyses were carried out. We show that clustering into two classes allows the conventional convective and stratiform classes to be determined, whereas classification into five classes allows this convective / stratiform classification to be further refined. Finally, our study made it possible to reveal the presence of some specific relationships between these five classes and the microphysics of their associated rain events.

## 1. Introduction

The analysis of "precipitation events" or "rain events" can be used to obtain information concerning the characteristics of precipitation at a particular location, and for a specific application. This is a convenient way to summarize precipitation time series in the form of a small number of characteristics that make sense for particular applications.

The concept of a precipitation event is not new, and has been used for many years (Eagleson, 1970; Brown et al., 1984). A wide variety of definitions, varying according to the context of each study, has been reported in the literature (Larsen and Teves, 2015). Moreover, when a rain rate time series (generally based on rain gauge records) is broken down into individual rainfall events, a wide variety of their characteristics, such as average rainfall rate, rain event duration and rainfall event distribution (known as hydrological information), can be computed for each event. Our analysis of the literature has led to the identification of seventeen features used to characterize rainfall, which makes it quite difficult to compare different studies. The first goal of the present study is to select a reduced set of features characterizing rainfall events, through the use of a data-driven approach, without taking *a priori* knowledge of the field of application into account, thereby characterizing rainfall events in the most parsimonious and efficient manner.

The second goal is to assess, without using any *a priori* criteria, whether the rain events are still correctly clustered by the most relevant observed features. Indeed, atmospheric process specialists distinguish between stratiform and convective events, arguing that the physical processes involved in their evolution are different. The goal here is to check that a small sample of variables, derived from spot measurements to describe rain events, can allow this distinction to be made, and ultimately be used to refine it. Hydrological (hereafter referred to as "macrophysical") information makes use of rain gauge measurements to characterize rain events. This information is defined in order to characterize the features of global events, but not to provide any information concerning the raindrop microphysics of the event. Nevertheless, in many applications such as remote sensing, knowledge of the microphysics is essential. One key parameter in remote sensing is the raindrop size distribution, noted as N(D), which is defined by the number of raindrops per unit volume and per unit raindrop diameter (D). Information related to the raindrop size distribution is often derived from its proxies, as explained in section 5. Such features are not currently accessible through rain-gauge measurements, which provide macrophysical information only. However, more expensive devices referred to as disdrometers can provide both hydrological and microphysical information. There are currently several tens of thousands of rain gauges operated throughout the world, in locations equipped with a far smaller number (if any) of disdrometers. As described later in this paper, it is possible to retrieve some microphysical information from the hydrological data. As a consequence, rain gauge data could provide valuable information in microphysics studies, through the use of a statistical approach to indirectly infer the missing microphysics information. In the following, the terms "macrophysical" or "hydrological" information are associated with characteristics related to rain rates or rain accumulation, whereas the term "microphysical" is associated with the characteristics of the raindrop size distribution.

In the present study, we use a data-driven approach to study the relationships between different rain properties. As disdrometers provide drop size distributions, they allow one-minute (or shorter) rain rates to be estimated, and in the present study these can be used to derive the hydrological information of interest, which is coherent with the data that would be provided by standard rain gauges. Through the combined interpretation of microphysical and hydrological information, we are also able to analyse the microphysical properties of the rain-event clusters provided by our algorithm. This makes it possible to retrieve (unobservable) microphysical information from rain gauge measurements.

From a single rain-rate times series, observed with a one minute time resolution, we seek to answer the following questions:

- among the large number of hydrological variables described in the literature, which are the most significant?

- does the resulting description of rain events allow different types of rain event to be discriminated?

- what (unobserved) microphysical properties of an event, or type of rain event, can be inferred from its macrophysical description?

Our paper is structured as follows. Section 2 presents the data used in our study, and lists various hydrological parameters that are commonly found in the literature. Seventeen macrophysical variables are identified, requiring appropriate normalisation. Section 3 presents our methodology, which is based on the use of a genetic algorithm (GA) implementing a self-organizing map (SOM also referred to as a topological map). This unsupervised approach is used to select a small subset of variables from the 17 identified variables, allowing a parsimonious characterisation of rainfall events to be applied. An exploratory statistical analysis of rainfall events is provided. In section 4, the rainfall events are grouped in clusters, and are divided into two classes. It is then shown that this grouping of the dataset corresponds to the standard convective / stratiform classification. We then propose a five-subclass classification, which corresponds to a refinement of the two initial groups. In section 5 we include some additional microphysical features of rainfall events, allowing the microphysical properties of the five previously defined event classes to be studied. Our conclusions are presented in section 6.

2

## 2. The disdrometer datasets - data processing methodology

This research relies on the analysis of raindrop measurements obtained with a Dual-Beam Spectropluviometer (DBS) disdrometer, first described by Delahaye et al. (2006). This instrument allows the arrival time, diameter and fall velocity of incoming drops to be recorded. As the capture area of the sensor is 100 cm$^2$ its observations can be considered, in spatial terms, to be "point-like". In the present study the integration time $T_{int}$ was set to one minute, and the raindrop measurements were used to estimate the corresponding one-minute rain-rate time series $RR_t(t)$. In order to eliminate false raindrop detections that could be generated by dust or insects, a threshold $T_0=0.1$ mm.h$^{-1}$ was applied. Rain rates lower than $T_0$ are thus set to zero. This conventional threshold is also chosen to ensure coherency with previous studies (Verrier et al., 2013; Llasat et al., 2001). In the present study, we worked with two datasets recorded during the period between July 2008 and July 2014, at the "Site Instrumental de Recherche par Télédétection Atmosphérique" (SIRTA[1]) in Palaiseau, France.

### 2.1 Rain event definition

In everyday life, it is common knowledge that rain starts at a certain moment, and stops some time later. However, due to its discreet nature, rain (which generally consists in a very large number of raindrops) is not an easy concept to define. Indeed, the exact definition of a rain event will depend on the sensor's characteristics (specific surface capture, detection threshold, instrumental noise), as well as the spatial or temporal resolution chosen for the study. This definition may also depend on the purpose of the study, and thus on the scientific community behind it. There is thus a wide range of criteria used to break down precipitation records into rain events. For this reason, it is important to define and apply an unambiguous definition of a "rain event".

In this study, the pattern produced by the one-minute rain rate time series $RR_t(t)$ can be simplified by grouping non-null rain rates into a set of separate "primitive events" (Brown et al., 1985). On the basis of an assigned Minimum Inter-event Time (MIT) (Coutinho et al., 2014) each rain rate value, corresponding to a specific one-minute period of observation, is assigned to a given rainfall event, i.e. either the rainfall event in progress, or a subsequent event that is considered to be independent and "new". The MIT could also be defined as the duration of a dry period $D_{dry}$ following which the next occurrence of non-null rainfall marks the beginning of a new event. For dry periods shorter than the MIT, rain rates from either side of this period are considered to belong to the same "composite event". Various authors have proposed different values of MIT that ensure event independence. Llasat (2001) noted that: "*The definition of an episode is quite subjective. In this case it was felt possible to distinguish between two different episodes, when the time which elapses between them without rainfall exceeds 1h, which ensures that, the two episodes come from different 'clouds'*". Bocquillon and Moussa (2014) wrote: "*the constant rain observations on less than thirty minutes represent only 5% of all the rainy periods. The representative threshold of the discretization of the data is 30 minutes to an hour.*"

Dunkerley (2008 a, b) carried out an analysis of the Inter-Event Time (IET) in order to check the influence of this variable on the definition of rainfall events, and its influence on the average rainfall rate. As emphasized in this study, when determining a value for the MIT, it is crucial to find an appropriate compromise between the independence of rain events and the intra-event variability of rain rates. The choice of MIT thus has a direct impact on the macrophysical characteristics that are ultimately determined by the analysis. Other researchers have proposed to use MIT values of 20 minutes, 1hour or even 1day (see Dunkerley, 2008a for a detailed list). In the present study it was decided to set the MIT to 30 minutes. This is in agreement with the value used by Coutinho et al. (2014), Haile et al. (2011), Dunkerley (2008a, b), Balme et al. (2006) and Cosgrove and Garstang (1995).

---

[1]  http://sirta-dev.ipsl.jussieu.fr/joomla/index.php/85-article-sans-categorie/71-sirta-home-page

When applied to our dataset, this choice leads to the identification of 545 rain events, which can be divided up into two subsets, i.e. one for learning and the other for testing (Tab. 1). The learning dataset is composed of observations collected over a two-year period between 2013 and 2014, with an availability of 96.4%, whereas the test dataset collected during the 2008 – 2012 period contains periods with missing data, due to a malfunction of the recording device.

**Table 1:  Observation periods and availability of DBS observations, and numbers of rain events for the learning and test datasets**

## 2.2 Macrophysical description of rain events

Rain events contain a wealth of information, which generally needs to be condensed into a limited set of well-chosen features. However, there is no conventional or commonly accepted list, or specific set of macrophysical features that can be used to accurately describe and summarize an event. In the present study it was thus decided to consider a large number of features, allowing the macrophysical rain event information described in the literature to be correctly represented. Seventeen characteristics were selected and identified (Llasat, 2001; Moussa, 1991), and are listed in Tab. 2. Some of these are parameter-dependent, such as $P_c$ which uses 3 values of the parameter c. These 3 values lead to 3 $Pc$ indices, namely $Pc_1$, $Pc_2$, $Pc_3$. Finally, a total of 23 descriptors were defined and numbered from 1 to 23 (column 1 in Tab. 2).

**Table  2: The 23 variables identified in the literature, used for the characterization of rain events**

Among the 23 indicators (hereafter referred to as variables) corresponding to the previously defined features, some are very well known. These include the event duration ($D_e$), the quartile ($Q_i$), the mean event rain rate ($R_m$) and the standard rain rate deviation ($\sigma_R$). Other less traditional parameters such as the parameter $\beta_L$ (indicator for the convective nature of the rain, see Llasat (2001)), the absolute rain rate variation of order c ($P_c$), or the absolute rain rate variation ($P_{s,c}$). Some variables that are usually used to describe time series, such as the fractal dimension, multi-fractal parameters, trend, seasonality, and autocorrelation, require a long series of data and are not well suited to an event-by-event analysis. This set of 17 features is not exhaustive, and some other features could also be included, depending on the application. One example is the case of hydrology, for which the positions of the intensity peaks inside the event could be a relevant feature. Although, for events comprising a very small number of samples (very low value of variable $D_e$), the computation of some indicators ($\sigma_R$, $Q_i$) is questionable, in the present study the 23 variables were computed for each of the 545 rain events.

## 2.3 PCA analysis and normalization step

It is important to note that very few of these 23 variables are compatible with the probabilistic assumptions generally associated with exploratory statistical methods. They are often highly variable, with highly skewed distributions, and therefore do not have normal distributions. It is thus more difficult to make direct use of standard statistical methods with this data, as it may lead to misleading interpretations (Daumas, 1982). It is thus necessary to introduce an additional step in order to transform the original distributions into *quasi* normally distributed distributions. The most suitable type of normalizing transformation for each of these variables was selected empirically, by testing 7 different possible transformations (Tab. 3).  For each variable, the retained transformation is that leading to a distribution having the strongest similarity to a normal distribution, i.e. with a kurtosis close to 3, and a skewness close to 0. For each indicator, the selected transformation is provided in the last column of Tab.2.

**Table  3: Transformations used to normalize the variables listed in Table 2.**

Following the normalisation step, Principal Component Analysis (PCA) was carried out on the learning data set (cf. end of section 2.1 and Tab.1). It follows that the two principal axes contain 73% of the total information, whereas the first 5 principal axes are needed to represent 90% of the total information. The $IET_p$ variable (#7) is very well correlated with axis 5, whereas the other variables are not. This means that there is no linear relationship between $IET_p$ and the other variables. For this reason, this variable was not considered as a possible candidate, in the variable selection process, during the remainder of the study. The results obtained in Section 4.2 confirm that there is no relationship between this variable and the other 22 variables. The correlation circle on axes 1 & 2 (Fig. 1.a.) shows that among the 23 variables, 16 are well correlated with the axis (close to a unit circle) and are distributed in approximately 5 groups (hereafter referred to as PCA groups). A first PCA group ($G_1$) can be identified by the variables, which are grouped close to the first axis, and are well correlated with it. As an example, this is the case for the variables $\sigma_R$ (#9), $P_{C_N}$ (#17 − 18) and β (#21 to 23). A second PCA group ($G_2$) comprises the variables $R_{max}$ (#11) and $P_{C3}$ (#16), just above axis 1. The third PCA group ($G_3$) is formed by the variable $P_{C2}$ (#15) only. The fourth PCA group ($G_4$) comprises the variables $P_{c1}$ (#14) and $Ps,c$ (#20) and is well correlated with axis 2. The last PCA group ($G_5$) is formed by the variable $D_e$ (#1). The correlation circle on axis 1 & 3 (fig. 1.b.) shows that the variables $Q_1$ (#4), $Q_2$ (#5) and $M_0$ (#10) are quite well represented by these two axes. A similar remark can be made for variables $D_d$ (#3) and $\beta_{L1}$ (#21) on axes 1 & 4 (not shown).

Finally, PCA analysis clearly shows that, within each PCA group, many variables are highly inter-correlated, i.e. linearly dependant on each other. This means that several variables could be removed with no substantial loss of information. This leads to the following question: which variables can be removed in order to retain the most parsimonious subset of variables representative of the full dataset? The PCA extracts summary variables, which are a linear combination of original variables, but do not allow for the selection of variables. To answer to this question, we propose a method for the global selection of variables, which seeks to identify the relevant variables in a dataset. As it appears to be intuitively more advantageous to select variables with a physical sense, rather than using dimension reduction methods (e.g. PCA which is more suitable for the detection of linear relationships), the proposed method is based on the use of a genetic algorithm. The following section provides a brief introduction into the concept of genetic algorithms, and shows how they can be advantageously used for the selection of variables in the context of the present study.

**Figure 1: PCA on the learning data set based on the 23 variables described in Tab. 3. Left: correlation circle on axes 1 & 2. Right: correlation circle on axes 1 & 3. All of the variables are normalised according to the last column of Table 2.**

**3. Variable selection using a genetic algorithm**

Computer-assisted variable selection is important for several reasons. Indeed, the selection of a subset of variables in a high-dimensional space can improve the performance of the model or its statistical properties, but also provides more robust models and reduces their complexity. In practice, it is not generally possible to try all potential combinations of variables, and to select the best of these, as a consequence of the enormous computational cost associated with such an approach. Among the many different variable-selection techniques described in the literature (Guyon and Elisseeff, 2003), we chose to develop a model based on the use of genetic algorithms, to search for an optimal subset of variables. Genetic algorithms (GAs) (Holland, 1975) are stochastic optimization algorithms based on the mechanics of natural selection, and the genetics described by Charles Darwin. In our study, a chromosome is defined as a subset made up from our 23 variables. A first generation composed of a population of 60 potential chromosomes is arbitrarily chosen. The performance of each chromosome (i.e. for each corresponding subset of 60 variables) is evaluated through a fitness function f. This fitness function is defined in such a way

that the higher its value, the greater the fitness function's ability to represent the full dataset (of dimension 23), using the smallest possible number of variables. On the basis of the performance of these 60 chromosomes, we create a new generation of 60 potential-solution chromosomes, using classical evolutionary operators: selection, crossover and mutation. The performance of this new generation is then evaluated. This cycle is repeated until a predefined stop criterion is satisfied. The best chromosome from the current generation then provides the optimal subset of variables.

## 3.1 Methodology

We define by $x^k$ the chromosome number k:  $x^k = (x_1, x_2, \ldots, x_{23})$. $x^k$  is a binary vector in $\{0,1\}^{23}$ space such that each component has the following meaning:

$$\forall i \in \{1, \ldots, 23\}, \quad \begin{cases} x_i = 1 & The\ variable\ number\ i\ in\ Tab.\ 2\ column\ 1\ is\ selected \\ x_i = 0 & The\ variable\ number\ i\ in\ Tab.\ 2\ column\ 1\ is\ not\ selected \end{cases} \qquad (1)$$

The word "selected" in eq. 1 means that the corresponding variable will be used, both in the learning step described in "Step 2" below, and for performance evaluation. Otherwise, if the corresponding variable is not selected it will be used only for performance evaluation.

As previously stated, the fitness function allows a measure to be provided of how well a minimal subset of variables can represent the entire data space (in dimension 23). The fitness function $f$ is thus defined as follows:

$$f(x^k) = \frac{1}{nb(x^k)\ te(x^k)} \qquad (2)$$

where

$x^k$ : chromosomes number $k$

$nb(x^k)$   :   number of selected variables in chromosome $x^k$

$te(x^k)$   :   topological error associated to chromosome $x^k$

As the aim of this approach is to minimise the number of selected variables $nb$ and the topological error $te$, we seek to maximize the fitness function. The estimation of the topological error made from a Self-Organizing Map (SOM) is somewhat complicated, and requires some explanation. The notion of a Self-Organizing Map, introduced by Teuvo Kohonen (Kohonen, 1982, 2001), makes use of a popular clustering and visualization algorithm. SOM is a neural network algorithm based on unsupervised learning, derived from the technique of competitive learning (Kohonen, 1982, 2001; Vesanto and Alhoniemi, 2000). It may be considered as a nonlinear generalization, which has many advantages over the conventional feature extraction techniques such as Empirical Orthogonal Functions (EOF), or Principal Component analysis PCA (e.g., Liu et al., 2006). SOM applications are becoming increasingly useful in geosciences (e.g., Liu and Weisberg, 2011). As stated by Uriarte and Martín (2008): "The SOM provides a nonlinear, ordered, smooth mapping of high-dimensional input data manifolds onto the elements of a regular, low-dimensional array. The main characteristic of the projection provided by this algorithm is the preservation of

neighborhood relationships; as far as possible, nearby data vectors in the input space are mapped onto neighboring locations in the output space". This property makes it straightforward to compute a topological error (see Uriarte and Martín, 2008, eq. 2). For each of the $x^k$ chromosomes, a Self-Organizing Map $M(x^k)$ is learned on the learning dataset. Only the selected variables are used during the learning process. Finally, for each Map $M(x^k)$, the topological error $te(x^k)$ can be computed in accordance with eq. 2 in Uriarte and Martín (2008). Section 4 provides additional information concerning Self-Organising Maps.

The Genetic Algorithm is based on the following five steps (Fig.2.):

*Step 1- Initialization:* (initial population)

Randomly generate a population $\{x^k, k=1, ..., 60\}$ of 60 chromosomes of dimension 23.

*Step 2- Evaluation:* For each of the $x^k$ chromosomes, a SOM $M(x^k)$ is learned. Only the selected variables are used for learning. Once the learning has been completed, the test dataset and the 23 variables are used on each of the 60 maps to estimate their topological error $te(x^k)$ allowing their fitness score $f(x^k)$ to be computed.

*Step 3-* Select the best chromosome $x^{Best}$ from the full set of 60 chromosomes according to the fitness score previously computed with the test dataset. If $x^{Best}$ remains unchanged over a period of 50 generations, stop the procedure and select the most relevant variables, i.e. those for which the corresponding components are equal to 1 in $x^{Best}$. Otherwise, go to step 4.

*Step 4- Selection:* Create a new population of 60 chromosomes from the current population, by randomly sampling with replacement chromosomes based on their probabilities, determined using the formula:

$$\Pr(x^k) = \frac{f(x^k)}{\sum_{i=1}^{60} f(x^i)} \tag{3}$$

*Step 5- Reproduction:* Mutation and Crossover possibilities in the new population.

Mutation: This consists in modifying (or not) certain components of the chromosomes. The probability of mutation is in general very low, and is commonly set to $p = 10^{-7}$. In the present case, the number of generations needed to reach the objective is less than a few hundred, such that the probability of a mutation is very low.

Crossover: In an initial step $\frac{60}{2} = 30$ pairs of chromosomes are randomly drawn from the population. Then, for each pair ($x^k$, $x^l$) (called parents) one crossover point, noted $I_c$, is randomly drawn over the range [1, 23], using a discrete uniform law. Two new chromosomes ($x^{k'}$, $x^{l'}$) are created as follows:

$$\begin{cases} x^{k'} = (x_1^k, x_2^k, ..., x_{I_c}^k, x_{I_c+1}^l, x_{I_c+2}^l, ..., x_{23}^l) \\ x^{l'} = (x_1^l, x_2^l, ..., x_{I_c}^l, x_{I_c+1}^k, x_{I_c+2}^k, ..., x_{23}^k) \end{cases} \tag{4}$$

Thus, from two parents, two children are generated, allowing a new generation to be produced with the same number of chromosomes. Finally, the algorithm returns to step 2.

**Figure 2: Diagram for the selection of variables based on a Genetic Algorithm associated with Kohonen Maps**

## 3. 2 Parsimonious description of a rain event

The genetic algorithm is applied to our datasets in order to obtain an optimal subset of variables forming a subspace, which can (in a certain sense) provide relatively accurate information concerning the global space, whilst having the particularity of containing non-redundant information. At the 187[th] generation the algorithm produces a subspace comprising 5 variables, namely : Event duration $D_e$ (#1), Standard deviation $\sigma_R$ (#9), maximum rain rate during event $R_{max}$ (#11), Rain event depth $R_d$ (#13), and Absolute rain rate variation $P_{c1}$ (#14).

The 3 variables ($D_e$, $R_{max}$, $R_d$) selected using this data-driven approach are commonly used in the study of hydrological processes (Haile and al., 2011). Moreover it should be noted that the commonly used variable $R_m$, which is computed simply by dividing the Rain event depth ($R_d$) by the duration ($D_e$), was not selected by the algorithm. This result could be expected, since it is correlated with the latter variables, and the algorithm provides a parsimonious description. Concerning the Absolute rain rate variation ($P_{c1}$), this variable was proposed by Moussa and Bocquillon (1991). It tends to provide information on the structure of the events, more specifically related to smooth events with a small number of sharp peaks. In fact, this variable promotes low variations of $RR_t$ because $P_{C1}$ is in a certain sense a structure function of order $c_1$ of the variable $RR_t$ (see #14 column 4 in Tab. 2), with a low value for the exponent ($c_1 = 0.5$). Finally, the standard deviation variable ($\sigma_R$), which is a second-order moment, is the most commonly used indicator to describe the variability of the precipitation rate within the rain event.

## 4. SOM learned with the five selected variables

A SOM is a topological map composed of neurons. In the present case, a neuron is a vector of dimension 23 containing the 23 variables defined in Tab. 2. Each neuron has 6 neighboring neurons. SOM is an unsupervised neural network trained by a competitive learning strategy that performs two tasks: vector quantization and vector projection. The SOM, which is different to k-means, uses the neighborhood interaction set to learn the topological structure hidden in the data. In addition, in order to achieve optimal referent vector (neuron) matching, its neighbors on the map are updated, leading to the generation of regions in which neurons located in the same neighborhood are very similar. The SOM can thus be considered as an algorithm that maps a high-dimensional data space onto a two-dimensional space called a map. A map can be used both to reduce the amount data by means of clustering, and to project the data in a nonlinear manner onto a regular grid (the map grid).

In the present study we used the toolbox developed by "the SOM Toolbox Team", which is available at the following site: http://www.cis.hut.fi/somtoolbox/. A SOM with 8×8 = 64 neurons is considered here. This choice corresponds to a compromise, since a smaller map would not be able to distinguish fine details whereas, in view of the number of observations, and a larger map would not be meaningful.

After learning by the GA algorithm described in the previous section, the resulting map $M(\boldsymbol{x}^{Best})$ can be used to assign to any event the best matching reference vector (neuron), in accordance with the 5 selected variables associated with the chromosome $\boldsymbol{x}^{Best}$. The $M(\boldsymbol{x}^{Best})$ map obtained with this procedure can be considered as an optimal representation of the initial data set.

Fig. 3 shows the distance matrix. For each neuron, the color indicates the mean distance between a neuron and its neighbors. The value at the center of each neuron represents the number of rain events of the learning data set, captured by the corresponding neuron. All neurons capture rain events and slightly more than half of these capture between 3 and 5 rain events, which is close to the value that would be obtained ($234 / 64 \cong 4$) if the rain events were uniformly distributed over the map.

8

**Figure 3: Distance matrix for the $M(x^{Best})$ map: The colour of each neuron represents the average distance between itself and its neighbouring neurons. The value inside each neuron indicates the number of rain events that it has captured inside the learning dataset. The black line separates the neurons into 2 classes, using the Hierarchical Ascendant Classification (see section 4.1). The arrows represent the gradients of the variables $R_m$, $\sigma_R$ and $D_e$**

## 4.1 Projection of the selected and unlearned variables onto the SOM

The five variables $D_e$, $\sigma_R$, $R_{max}$, $R_d$ and $P_{c1}$ used for learning are referred to as 'selected' variables, whereas the remaining 18 variables are referred to as 'unlearned' variables. In order to study the relationship between these variables, Fig. 4 shows the projections for each of the variables in the $M(x^{Best})$ map obtained with the aforementioned GA selection algorithm. The variables are discussed individually, by considering their structure, as well as the relationships between them. We note that the map is well structured for the majority of variables. This advantageous structuration of most of the variables confirms the ability of the selected variables to summarize all of the significant characteristics of rain events. Only a small number of characteristics are not adequately represented. It should be noted that almost all variables are structured according to the first or second diagonal. Among these, one may consider an initial subset comprising variables that are more or less structured according to the first diagonal. This is the case for the unlearned variable $D_d$, as well as for the selected variables $P_{c1}$ and $D_e$. A second subset comprising variables that are structured in approximate accordance with the second diagonal can be identified. This is the case for the unlearned variables $R_{m,r}$, $R_{m,}$, $P_{C_{Ni}}$ which are very similar to the selected variable $\sigma_R$. The unlearned variables $Q_3$, $P_{C3}$, $P_{S,C}$ also belong to the second subset and have a structure close to that of the selected variable $R_{max}$.

The map can be related to the previously implemented PCA (Fig.1). As can be seen in Fig. 4, the variables $P_{C3}$ (#16) and $R_{max}$ (#11), which have a similar structure, also belong to the same PCA group, namely group $G_1$ (see section 2.3, Fig. 1.a). It is interesting to note that the variables $P_{S,C}$ (#20) and $R_{max}$ (#11), which also have a similar structure, do not belong to the same PCA group (groups $G_4$ and $G_2$ respectively) and are uncorrelated (they are orthogonal in Fig.1a). This remark means that the topological map reveals a relationship that cannot be detected using PCA. As the Rain event depth ($R_d$) depends on both the duration and the intensity of the events, the corresponding map has a top-down structure. Two distinct situations thus occur:

- Those events which contribute the greatest quantities of water (Fig. 4., brown neuron at the bottom right of $R_d$) are among the longest (see corresponding neuron of $D_e$), but do not have an extremely high peak rain rate (see corresponding neuron of $R_{max}$) and are quite smooth (see corresponding neuron of $P_{c1}$ and $\sigma_R$).

- Other events which contribute large amounts of water (but less than previously) (Fig. 4. red neuron at the bottom left of $R_d$) have short durations (see corresponding neuron of $D_e$), but are violent (see corresponding neuron of $R_{max}$) and are less smooth (see corresponding neuron of $P_{c1}$ and $\sigma_R$). The latter case reflects situations that are typical of convective storms.

The resulting map confirms the dependence structure of the two hydrological variables $R_d$ and $D_e$ studied by Gargouri and Chebchoub (2010).

The variable $IET_p$ (Previous IET): the map is not structured, reflecting the independence of the characteristics of a rain event with respect to the drought period preceding the event. This corroborates the results of several previous studies (Lavergnat and Gole, 1998, 2006; Akrour et al., 2015) dealing with rain support simulations. When studying temperate mid-latitudes for relatively short periods, these authors noticed that successive rain and no-rain periods are uncorrelated, such that a rain time series could be considered as an independently drawn, alternating series of rain events and periods without rain. This is equivalent to an inter-event time (IET) that does not characterize the rain events. The same effects are not necessarily observed at other locations, and under different climatological conditions. Brown et al. (1983) also investigated a possible correlation between IETp and the intra-event characteristics, and concluded that their data provided no evidence of this.

Since the variable $\beta_L$: $\beta_L$ (Llasat, 2001) is considered to represent a measure of the convective nature of the rain, it makes sense that the three variables $\beta_{L1}$, $\beta_{L2}$, $\beta_{L3}$ are structured similarly, with the peak rain rate variable $R_{max}$. This relationship is clearly visible on the maps.

Several other relationships, which are not described in detail here, can be observed. These include the correlation between the Normalized Absolute rain rate variation ($P_{C_{Ni}}$) and the standard deviation of the intensity ($\sigma_R$). We conclude that the combination of the five selected variables provides a relatively accurate summary of the information needed to describe the rain events. The poor structuring of some variables is justified by the independence of these variables with respect to the properties of the rain events; this is the case for the variable Dry Percentage in event $D_d$ or the variable $IET_P$.

Figure 4: Projection of the $M(x^{Best})$ map according to the 23 variables. The red-framed variables are those selected by the GA algorithm. The last two variables $D_m$ and $N_0^*$ are defined in section 5

**4.2 Representation of rain events on SOM**

In an effort to provide additional information for validation of the map, we compared each of the 23 variables with their corresponding value given by the SOM, for the learning dataset and the test dataset. For each of the 311 events of the test dataset, the best matching unit of the SOM, i.e. the neuron that is the closest to the event, is determined with respect to the five selected variables. As an example, for each event Fig. 5 shows the current value of the unlearned variable $\beta_{L3}$ as a function of the corresponding value given by the best matching unit of the event. A spread can be seen, in particular in the central zone, whereas the spread is relatively small for values located near to the edges (which are more numerous). A linear regression leads to a relatively good determination coefficient (R-square) (0.96 and 0.89 respectively, for the learning and test data sets). Table 4 lists the value of R-square for the 23 variables obtained with the learning and test data sets. As expected, the coefficient of determination of the variable I$ET_p$ is very poor (*0.31/0.26*), since this variable is not related to the 5 selected variables, and as a consequence cannot be well represented by the SOM (Fig. 4). The selected variables have good determination coefficients, with both the learning and the test datasets; this confirms the quality of the learning and the generalization ability of the SOM. The quality of the learning step is confirmed by the fact that the R-square values of the selected variables obtained on the test set are close to those obtained on the learning set. The R-squares corresponding to the unlearned variables obtained on the learning data set emphasize the ability of the selected variables to provide the information contained in the unlearned variables; in the case of the test data set it denotes the ability of the SOM to derive all event characteristics from the selected variables only.

Figure 5: the variable $\beta_{L3}$ versus its corresponding value, given by the best matching unit: from the learning data set (circles), and on the test dataset (stars). The solid line corresponds to the first diagonal

Table 4: Coefficient of determination obtained on the learning and test data sets. The values with a dark grey background correspond to the 5 selected variables

**4.3 Hierarchical clustering of rain events**

We have shown that the distance matrix (Fig. 3) confirms the successful deployment of the map. Based on the distance between neurons, it appears that neurons can be grouped to obtain a limited number of classes, each with its own characteristics. In order to group the 64 neurons into a small number of classes, a hierarchical cluster analysis was carried out (Everitt, 1974).

Only the five selected variables were used for the classification, and a Euclidian distance was selected for the hierarchical algorithm. Fig. 6 shows the resulting dendrogram, applied to the 64 neurons.

**Figure 6  Dendrogram obtained from the Hierarchical Cluster Analysis of the 64 neurons in the SOM. The horizontal dashed line represents the threshold between the two classes.**

Depending on the physical processes involved, experts tend to separate rain events into two different classes: stratiform and convective events. Although this classification is relatively crude, since stratiform and convective events can sometimes exist inside the same rain event, it is very commonly used. Concerning the times series, most authors use a very simple scheme to distinguish between stratiform and convective rain types. For reasons of simplicity, rain classification is sometimes defined using the instantaneous rain rate and the standard deviation estimated over consecutive samples. As an example, Bringi et al. (2003) defined stratiform rain samples when the standard deviation of the rain rate, taken over five consecutive 2-min samples, is less than 1.5 mm.h$^{-1}$, the convective rain samples are defined for a rain rate greater than or equal to 5 mm.h$^{-1}$, and the standard deviation of the rain rate over five consecutive 2-min samples is greater than 1.5 mm.h$^{-1}$.

Firstly, we separate the dendrogram into two classes. The first class contains 51 neurons and 79% of the observations, whereas the second class contains 13 neurons and 21% of the observations. The solid black line in Fig. 3 corresponds to the dividing line between these two classes. The first class, containing the greatest number of neurons, is in most cases characterized by relatively low rain rates. This can be seen by examining the structure of the map, according to the mean rain rate variable ($R_m$). Moreover, analysis of the standard deviation (small values of $\sigma_R$), absolute rain rates $P_c$ (high values of $P_{c1}$, and low values of $P_{c3}$) shows that this class is more or less characterised by quiet, homogeneous events. Our analysis of event durations ($D_e$) shows that this class contains both short and long durations, but is dominated by the latter. These characterizations are relatively well matched to a description involving stratiform and stable precipitations, which are often the consequence of the slow, large-scale uprising of a large mass of moist air which then condenses uniformly.

The second group is characterized by a smaller number of neurons. This corresponds to the higher values of the mean rain rates ($R_m$) and peak rain rates ($R_{max}$). The variables $\sigma_R$, $P_c$ have the opposite values with respect to those of the previous group. Most of the event durations ($D_e$) in this group are short, with the exception of neuron #64 (bottom right on the maps). This group fits well with the definition of convective events resulting from the rapid rise of air masses loaded with moisture, for buoyancy. This convective moist air can lead to the development of cumulus clouds up to an altitude in excess of 10 km, and to heavy rain.

Our analysis of the structure of the variables $\beta_{L1}, \beta_{L2}, \beta_{L3}$ in Fig. 4 confirms the previous interpretation of the two groups. These three variables, which are representative of convective rain, have high values for the neurons belonging to this group.

Figs. 7.a and 7.b show the neurons in the $R_m$, $\beta_{L3}$ and $P_{C2}$ subspace. These 3 variables were not used in the learning step. Nevertheless, the two classes are well separated, although an overlap does occur in Fig 7a due to neuron #64 (bottom right on the map, Fig. 4). We checked although it belongs to the convective class, this neuron nevertheless has some characteristics of the stratiform class.

**Figure 7: Representation of the neurons in the Rm, $\beta_{L3}$ and Rm, and P$_{C2}$ subspaces. The stars represent neurons from group 1 (stratiform), and the squares correspond to neurons from group 2 (convective). Dashed lines indicate the neuron #64**

The hypothesis that the two categories of precipitation events corresponding to different dynamic regimes can be identified solely on the basis of hydrometeorological variables is in agreement with the findings of Molini et al. (2011). These authors

11

have shown that there is a strong agreement between the hydro-meteorological classification (based on the duration and extent of events from rain gauge network data), and dynamic classifications (the convective adjustment time-scale identified to distinguish between equilibrium and non-equilibrium convection derived from ECMWF analysis). We conclude that this unsupervised automatic clustering, based on the five selected variables, makes it possible to correctly implement a classification with these two well-known classes (stratiform and convective). It should be noted that, unlike other classifications described in the literature, this was established without making use of *a priori* information, since it is produced by an unsupervised process.

## 4.4 Classification of events into several classes

From the stratiform and convective classification described above, it is interesting to refine the two classes into a set of subclasses. The synoptic rainfall associated with mid-latitude depressions provides an example of stratiform precipitation, which forms in depressions in the vicinity of warm and cold fronts. The very light type of rainfall (drizzle) associated with stratus or stratocumulus is included in the class of stratiform precipitation. This can occur under anticyclonic conditions, or in the warm region of a depression. The associated rain depths ($R_d$) are minimal, and usually have no hydrological impact other than superficial wetting. In order to identify relevant subclasses, our classification was broken down into a number of unknown subclasses, such that $n > 2$.

An important step in hierarchical clustering is the selection of an optimal number of partitions ($n_{opt}$) in the dataset (Grazioli et al., 2015). Many indices can be used to evaluate each partition, from the point of view of data similarity only. Most of these evaluate the scattering inside each cluster, with respect to the distance between clusters, and assign relatively favourable scores to partitions with compact and well-separated clusters. Although different indices were tested, these did not provide the same number of subclasses (between 2 and 32 with the indices tested in this study). It should be noted that these did not take the physical meaning of each class into account. Finally, we chose $n_{opt} = 5$, since higher values led to classes with the same physical sense. The new classification based on the use of five subclasses is shown in Fig. 8.

**Figure 8: Hierarchical Clustering of the map into five subclasses. The colours represent the subclass numbers: Subclass 1: Dark blue, Subclass 2 : blue, Subclass 3 : Green, Subclass 4 : Orange, Subclass 5 : Red**

From these five subclasses, two belong to the stratiform class and the other three belong to the convective class. In the learning dataset, the first subclass represents 12% of all events, and respectively 68%, 1.2%, 6.8% and 12% for subclasses 2 to 5. The characteristics of these five subclasses are summarized below and in Tab. 5. The five selected variables are remarkably heterogeneous between classes, meaning the accuracy of these variables for clustering*:*

**- Subclass 1** (drizzle and very light rain):  the main feature of this class is the very low mean value ($R_m$) and standard deviation $\sigma_R$ of rain rate events, in addition to the features of the superclass. The mean rain rate events lie in the range [0, 0.5] mm.h$^{-1}$ with a mean value of 0.36 mm.h$^{-1}$ and $\sigma_R$ in the range [0, 3] mm.h$^{-1}$, with a mean value equal to 0.1 mm.h$^{-1}$. Although this event has a significant duration, the corresponding subclass, which corresponds to drizzle, involves only small quantities of water. It can also be noted that a low value of $\beta_{L3}$ is a good indicator (< 0.01) for drizzle.

**- Subclass 2** ("normal" events): this is a relatively broad class containing 68 % of all events, with a mean event rain rate ($R_m$) in the range [0.5 , 6] mm.h$^{-1}$ and a mean value of 1.48 mm.h$^{-1}$. The standard deviation $\sigma_R$ lies in the range [1, 10] mm.h$^{-1}$ with a mean value of 2. This subclass is characterised by a significant relative variation of some parameters ($D_e$, $R_m$, $P_{c1}$ for instance), together with dry periods ($D_d$), which may be sufficiently long.

The three remaining subclasses correspond to convective classes of events, which are characterized by a strong temporal heterogeneity and significant intensities. Depending on the depth of rain events, this convective class is subdivided into:

**- Subclass 3:** containing relatively long events ($D_e$) with high values for the rain event depth ($R_d$) variable and $P_{c1}$. This class represents events with a very small likelihood of occurrence (1.2%).

**- Subclass 4:** containing relatively short events ($D_e$) with peak rain rate $R_{max} > 50$ mm.h$^{-1}$, in addition to strong heterogeneities ($\sigma_R$, $P_{C2}$ and $P_{C3}$ are high) and large values for the convective indicator ($\beta_{L3}$).

**- Subclass 5:** the events in this subclass are characterised by relatively low values for the rain event depth ($R_d$). This is due to the short duration of the events ($D_e$). The variables $\sigma_R$ and $P_{C3}$ remain high. Another feature of this subclass is that it includes continuous events only, with no short, embedded dry periods (low values of $D_d$ in Fig. 4 and Tab. 5).

**Table 5: Summary of rain event subclasses computed with the learning dataset.**

To conclude this section, this new classification allows the conventional definition for stratiform events to be refined. The convective classification can be subdivided into five different subclasses, each of which is homogeneous. This classification is obtained for mid-latitude climates. As the dataset used in this study is representative of only one specific region and topography (i.e. the temperate climate encountered in the Ile de France region, France), its analysis cannot reveal information related to different processes, i.e. those which are not sampled in the dataset. Such processes could lead to the identification of additional specific clusters of events. In particular, there are no orographic rainfall events or oceanic observations. The final step in this study involves assessing whether the homogeneous character of each class is preserved at the microphysics scale, and attempting to identify any relationships between the information present at the scale of both the microphysics and the macrophysics of these events (hydrological information).

## 5. Microphysical point of view

Our study of the microphysical properties of rain is based on a comprehensive analysis of its drop size distribution *N(D)*, corresponding to the number of raindrops per unit volume and per interval of diameter *D*. The shape of *N(D)* reflects the microphysical processes involved. The identification of various features of the drop size distribution, as well as the type of precipitation, is very useful for many applications. As an example, this information is used in the calculation of heating profiles in the precipitation parameterization of atmospheric models, to gain a more detailed understanding of microphysical processes, as well as for the development of rain retrieval algorithms applied to remote sensing observations. The microphysical characteristics of rainfall act as hidden variables that affect the relationship between microwave remote sensing measurements and the volume of water in a rainfall event (Ulaby, 1981; Iguchi, 2009). It can thus be very useful to use conventional rain gauges to determine the microphysical characteristics of rainfall events, thereby improving the quality of active or passive remote sensing observations, and the spatial properties of rainfall events in particular.

A general expression for the drop size distribution defined by Testud et al. (2001) is commonly used in the literature. This allows a distinction to be made between the stable shape function *f* and the variability induced by rain. This variability is represented by two microphysical parameters, namely the mass-weighted volume diameter ($D_m$) and the parameter $N_0^*$. In some studies, the term $N_w$ is used rather than $N_0^*$. Not all authors use exactly the same units, in particular Bringi et al. (2003) and Suh et al. (2016) use mm$^{-1}$m$^{-3}$ for the units of $N_w$, rather than the unit m$^{-4}$ which is used in this study for $N_0^*$.

$$N(D) = N_0^* f\left(\frac{D}{D_m}\right) \qquad [\text{m}^{-4}] \qquad (5)$$

where $D_m$ and $N_0^*$ are defined as:

$$D_m = \frac{M_4}{M_3} \; [\text{mm}] \;, \quad N_0^* = \frac{4^4}{\Gamma(4)} \frac{M_3^5}{M_4^4} \; [\text{m}^{-4}] \qquad\qquad (6)$$

and $M_i$ is $i$th-order moment of the drop size distribution $N(D)$:

$$M_i = \int_0^{+\infty} N(D) D^i dD \qquad\qquad (7)$$

Rain samples are usually analysed by computing the microphysical parameters ($D_m$ and $N_0^*$) for each rain sample obtained over a given time scale. In the present study, $N(D)$ is obtained by considering the entire raindrop collection corresponding to each rain event of (variable) duration $D_e$. This approach leads to one pair ($D_m, N_0^*$) of microphysics variables per rain event, whereas most other authors rely on values computed over a fixed time scale.

Projections of the learned map, according to $D_m$ and $N_0^*$, are shown in Fig. 4 (bottom right). It can be seen that the two maps are well structured, and that these two parameters have opposite influences on the map projection. Although these two microphysical parameters were not learned, the relationship between them is clearly accounted for by the information used to structure the map (the 5 selected variables). Moreover, the existence of a relationship between the microphysical and macrophysical features of the rainfall is also confirmed in this figure, since both of the macrophysical variables used to learn the SOM, i.e. $\sigma_R$ and $R_{max}$, have patterns similar to those revealed on the $D_m$ map.

Many authors, including Atlas et al., 1999; Bringi et al., 2003; Marzuki et al., 2013; Suh et al. 2016, have endeavoured to associate specific microphysical properties with each type of precipitation (convective or stratiform). In view of the maps shown in Fig. 4 and the convective/stratiform classification developed in section 4.3, we are able to confirm that precipitation events classified as stratiform express small values for $D_m$ and large values for $N_0^*$. In the case of the convective class, the opposite trend is observed (i.e. larger values for $D_m$ and smaller values for $N_0^*$). Similar observations have been reported by Testud et al. (2001). It can also be noticed that the two microphysical variables are relatively homogeneous in the convective class, whereas in the stratiform class they are characterised by a higher level of variability.

In order to improve our analysis of the microphysical information embedded in the dataset, we analysed the relationship between the two microphysical parameters using the reference vectors (neurons) from the map, which include information related to the original rain events.

Fig. 9 shows the variable $D_m$ as a function of $N_0^*$ for the 64 neurons on the map. This relationship is indicated through the use of distinct markers to identify the five subclasses defined in section 4.4, thus facilitating the discussion of the microphysics associated with stratiform and convective rain. The two solid lines show the linear regressions computed for these two classes.

In the case of the stratiform subclasses (1 and 2) a clear relationship can be observed between the two variables. The microphysics characteristics of these two subclasses are clearly distinct. Indeed, subclass 1 (drizzle and light rain) has the smallest $D_m$ and the highest $N_0^*$, and varies over just a small range. Conversely, as in the case of the macrophysical variables (see section 4.4), the microphysical characteristics of subclass 2 (normal events) are considerably more heterogeneous. Knowledge of $D_m$ makes it straightforward to identify the corresponding subclass. As a consequence, an event with $D_m$ lying in the range [0.5, 1] millimetre belongs to subclass 1. Similarly, it is very likely that an event with $D_m$ lying in the range [1, 1.7] millimetre belongs to subclass 2.

For the convective events (subclasses 3, 4, 5), small differences can be noticed with respect to $N_0^*$. In the range [1.7 , 2.5] mm, two neurons belonging to subclass 4 are close to a neuron belonging to subclass 5, and therefore have similar microphysics. Although they are located far from all other subclass 2 neurons, three isolated neurons belonging to subclass 2 (stratiform) can be noted. These are characterised by relatively strong values of $D_m$ (2 mm) and low values of $N_0^*$. The corresponding events

are a mixture of stratiform and convective rain. A typical case is given by convective rain associated with strong rain rates occurring at the beginning of an event, whereas the remainder of the event is stratiform with low rain rates and small variations.

*Figure 9: Microphysical variable $N_0^*$ versus $D_m$ for the five rainy event subclasses. The three neurons corresponding to mixed events are circled. The dashed lines indicate the limits defined by $D_m > 1.66$ and $Log(N_0^*) > 6.15$*

Following our classification, Fig. 9 indicates that there are real relationships between the macrophysical and microphysical variables. Nevertheless, knowledge of the variables $(D_m, N_0^*)$ does not allow the correct subclass to be determined in all cases.

Researchers who study microphysical features and their association with specific types of precipitation use simple schemes, based on rain rate estimations over a fixed period of integration (a few minutes), in order to separate stratiform and convective rain types. They also use these simple schemes to label $D_m$ and $N_0^*$ as stratiform or convective (Testud et al., 2001). This approach is significantly different to the method presented here, which assumes that all of the samples in a given event belong to the same class. Our values for $N_0^*$ and $D_m$ are thus computed for the time scale of a given event, rather than for a fixed integration time. Thus, although in the present study a good agreement is found for the range of values covered by $D_m$, those determined for $N_0^*$ do not cover the same range as in the case of the previously cited studies.

Many previous authors have observed that the drop size distribution is closely related to processes controlling rainfall development mechanisms. In the case of stratiform rainfall, the residence time of the drops is relatively long, and the raindrops grow by the accretion mechanism. In convective rainfall, raindrops grow by the collision–coalescence mechanism, associated with relatively strong vertical wind speeds. Numerous studies have been published concerning the variability of $N_0^*$ and $D_m$: Bringi et al. (2003) studied rain samples from diverse climates and analysed their variability in stratiform and convective rainfall; Marzuki et al. (2013) investigated the variability of the raindrop size distribution through a network of Parsivel disdrometers in Indonesia; and Suh et al. (2016) investigated the raindrop size distribution in Korea using a POSS disdrometer. In the case of stratiform rain, all of these authors observe that $N_0^*$ and $D_m$ are nearly log-linearly related, with a negative slope. This is consistent with the trend shown in Fig. 9 for the two stratiform subclasses (1 and 2). Even the three distinct neurons, which are isolated from the others, appear to be governed by the same relationship.

Marzuki et al. (2013) noted that during convective rain the increase in value of $N_0^*$ with decreasing $D_m$ is nearly log-linear, with a flatter slope. In the present case, the dependence is also log-linear, with a slope that is slightly flatter for convective events than for stratiform events. In the aforementioned studies, the data was aggregated over time, campaign or site, on the basis of a criterion computed over a fixed period of time. We believe that this process is weakly suited to determining the properties of convective events, as a consequence of their strong variability and shorter characteristic time. In this study we were able to retrieve the log-linear relationship between $N_0^*$ and $D_m$ without having to learn it directly.

When applying our algorithm to the various macroscopic properties by rain event, we also take into account the variability of rain within an individual rain event. Fig. 9 clearly shows that the spreading of parameters $N_0^*$ and $D_m$ inside each subclass has the same magnitude as the distance between subclasses. This remark confirms the hypothesis of Tapiador et al. (2010): the intra-event variability can exceed the inter-event variability, due to events arising from different precipitation systems. It is thus preferable to examine the properties of events with a more general approach, rather than using individual samples to study the distinction between stratiform and convective processes. The three isolated neurons in subclass 2 described above (circled in Fig. 9) have the same properties as the other events of their subclass (i.e. the same slope for the log-linear relationship between $N_0^*$ and $D_m$). This example confirms the ability of our methodology to preserve the macroscopic information needed to cluster rain events, thus allowing the intra-event variability as well as microphysical information to be (partially) retrieved. Suh et al. (2016) also compare $\log(N_0^*)$ and $D_m$ pdf, for the case of stratiform and convective samples over a 4-year period. On the basis of the $D_m$ pdf of both stratiform and convective classes, they compute a threshold value for $D_m$, such that when $D_m >$

1.66 mm the rainfall samples are mainly convective, and when $D_m < 1.66$ mm they are mainly stratiform. This finding is consistent with the results of Atlas et al. (1999), who also found a threshold value for $D_m$, distinguishing between convective and stratiform rainfall. In Fig. 9, it can be seen that this threshold is confirmed (vertical solid line), with $D_m$ smaller than 1.6 mm corresponding to stratiform events, whereas higher values correspond to mainly convective events. When we consider the events analysed in the present study, there are also three neurons corresponding to a "mixed event" beyond this threshold.

Suh et al. (2016) show in Fig. 4c of their study that the pdf for convective rainfall is higher than that corresponding to stratiform rainfall, when $\log(N_0^*) > 6.2$ ($N_w = 3.2$ in their figure). As described above, by considering the data corresponding to rain events, rather than to samples recorded over fixed periods of time, our range of values for $N_0^*$ is smaller than that used in other publications. In addition, $\log(N_0^*) < 6.15$ for all neurons labelled as convective in our study, which is very close to the value of 6.2 determined by Suh et al. (2016).

In view of the generally satisfactory retrieval of microphysical information from macrophysical parameters, we are of the opinion that the topological map successfully restores some of the information implicitly embedded in the dataset. It is thus interesting to note that the macrophysical parameters of rainfall are related to its microphysical properties. Firstly, the map collects similar events, whilst ensuring, through the minimization of topological errors, that the unfolding of the map is correct. A neuron is thus closer to its neighbours than to any other neuron on the map. This criterion ensures that the data space is optimally partitioned into connected subparts, such that the neurons on the map can be related to the underlying processes governing rainfall.

### 6. Conclusion

Although the definition of a 'rain xc event' is relatively subjective, this study underlines the advantages of using event analysis rather than sample analysis. This data-driven analysis of events shows that rain events exhibit coherent features. As a consequence of the discrete and intermittent nature of rainfall, some of the features commonly used to describe rain processes are inadequate, in particular when they defined for a fixed duration. Excessively long integration times (hours or days) can lead to the mixing of observations that correspond to distinct physical processes, and also to the mixing of rainy and clear air periods, within the same sample. An excessively short integration time (seconds, minutes) leads to noisy data, which is sensitive to the sensor's characteristics (sensor area, detection threshold and noise). By analyzing entire rain events, rather than short individual samples of fixed duration, it is possible to clearly identify certain relationships between the different features of rain events, in particular the influence of the microphysical properties of rain on its macrophysical characteristics. This approach allows the intra-event variability caused by measurement uncertainties to be reduced, thus improving the accuracy with which physical processes can be identified.

Once an event has been clearly identified, it is possible to choose a small number of variables to describe it. We present a new data-driven approach, which can be used to select the most relevant variables for this characterization. This approach has generic properties and can be adapted to many multivariate applications. A genetic algorithm, when combined with Self-Organizing Map (SOM) clustering, can allow the unsupervised selection of an optimal subset of five macrophysical variables. This is achieved by minimizing a score function, which depends on the topology error of the SOM and the number of variables. This score provides a parsimonious description of the event, whilst preserving as much as possible the topology of the initial space.

Numerous variables derived mainly from rain rate recordings are used to describe precipitation in the context of rain time series studies, and a wide variety of topics of interest, including hydrology, meteorology, climate, and weather forecasting. The algorithm proposed in this study produces a subspace formed by only 5 of the 23 rain features described in the literature. We show that these five features can be selected by the algorithm in an unsupervised manner and, from the macrophysical point of view, can provide an adequate description of the main characteristics of rainfall events. These characteristics are: the

event duration, the peak rain rate, the rain event depth, the standard deviation of the event rain rate, and the absolute rain rate variation of order 0.5.

In order to confirm the relevance of the five selected features, we analyze the corresponding SOM and are able to clearly reveal the presence of relationships between these features. This approach also reveals the independence of the inter-event time ($IET_p$) characteristic, and the weak dependence of the Dry percentage in event ($D_{d\%e}$) characteristic, thus confirming that a rain time series can be considered as an alternating series of independent rain events, interrupted by periods without rain. Hierarchical clustering allows the well-known separation between stratiform and convective events to be clearly identified. This dual classification is then refined into a set of five relatively homogeneous subclasses. The stratiform class is divided into 2 subclasses: a drizzle / very light rain subclass, and a normal event subclass. The convective class is divided into 3 subclasses, characterized by a strong temporal heterogeneity and significant rain rates.

As this research was based on the analysis of observations made in mid-latitude plains in France, the relevance of this classification remains to be confirmed through the analysis of datasets recorded in different climatic zones, and under different meteorological conditions, such as those encountered in mountainous or coastal areas. If the SOM described in the present study were learned with a more exhaustive dataset, a larger map would be produced, and this could reveal new types of rainfall behavior, which remained undetected in the current dataset. This point will be addressed in future studies.

The data-driven analysis of entire rain events (rather than the analysis of fixed-length samples) is relevant to the study of interactions between the macrophysical (based on the rain rate) and microphysical (based on raindrop) properties of rain. In the present study, several strong relationships were identified between these microphysical and macrophysical characteristics, and we show that some of the five subclasses identified in this analysis have specific microphysical characteristics. When a relationship between the microphysical and macrophysical properties of rain is identified, this can have many practical implications, especially for remote sensing. In the context of weather radar applications, the microphysical properties of rain are needed in order to estimate rain rates, through the use of the Z–R relationships. The estimation of microphysical rain characteristics, based on easily observable rain gauge measurements, could play a significant role in the development of the quantitative precipitation estimation (QPE).

## References

Akrour, N., Chazottes, A., Verrier, S., Mallet, C., and Barthes, L.: Simulation of yearly rainfall time series at microscale resolution with actual properties: Intermittency, scale invariance, and rainfall distribution. Water Resources Research, 51(9), 7417–7435, 2015.

Atlas, D., Ulbrich, C.W., Marks, F. D., Amitai, E., and Williams, C. R.: Systematic variation of drop size and radar-rainfall relations, J. Geophys. Res.-Atmos., 104, 6155–6169, 1999.

Balme, M., Vischel, T., Lebel, T., Peugeot, C., and Galle, S. : Assessing the water balance in the Sahel: impact of small scale rainfall variability on runoff Part 1: rainfall variability analysis. Journal of Hydrology 331:336–348, 2006.

Bringi, V. N., Chandrasekar, V., Hubbert, J., Gorgucci, E., Randeu, W. L., and Schoenhuber, M.: Raindrop size distribution in different climatic regimes from disdrometer and dual-polarized radar analysis, J. Atmos. Sci., 60, 354–365, 2003.

Brown, B.G., Katz, R. W., and Murphy, A.H.: Statistical analysis of climatological data to characterize erosion potential: 1. Precipitation Events in Western Oregon. Oregon Agricultural Experiment Station Spec. Rep. No. 689, Oregon State University (1983).

Brown, B.G., Katz, R. W., and Murphy, A.H.: Statistical analysis of climatological data to characterize erosion potential: 4. Freezing events in eastern Oregon/Washington. Oregon Agricultural Experiment Station Spec. Rep. No. 689, Oregon State University, 1984.

Brown, B.G., Katz, R. W., and Murphy, A.H.: Exploratory Analysis of Precipitation events with Implications for Stochastic Modeling. Journal of Climate and Applied meteorology(57-67), 1985.

Cosgrove, C.M. and Garstang, M.: Simulation of rain events from rain-gauge measurements. International Journal of Climatology 15, 1021–1029, 1995.

5  Coutinho, J.V., Almeida, C. Das, N., Leal. A.M. F., Barbarosa, L. R.: Characterization of sub-daily rainfall properties in three rain gauges located in northeast Brazil. Evolving Water Resources Systems: Understanding, Predicting and Managing Water–Society Interactions Proceedings of ICAR 2014, Bologna, Italy, 345-350, 2014.

Daumas, F. : Méthodes de normalisation de données, Revue de statistique appliquée, 30(4), 23-38, 1982.

Driscoll, E. D., Palhegyi, G. E., Strecker, E. W., & Shelley, P. E. (1989). Analysis of storm events characteristics for selected
10  rainfall gauges throughout the United States. *US Environmental Protection Agency, Washington, DC.*

Dunkerley, D.: Rain event properties in nature and in rainfall simulation experiments: a comparative review with recommendations for increasingly systematic study and reporting, Hydrological Processes, 22(22), 4415–4435, 2008a.

Dunkerley, D.: Identifying individual rain events from pluviograph records: a review with analysis of data from an Australian dryland site, Hydrological Processes, 22(26), 5024–5036, 2008.

15  Eagleson, P. S.: Dynamic Hydrology, McGraw-Hill, 1970.

Galmarini, S., Steyn, D. G., and Ainslie, B.: The scaling law relating world point-precipitation records to duration. Int. J. Climatol., 24, 533–546, 2004.

Everitt, B.: Cluster Analysis. London: Heinemann Educ. Books, 1974.

Delahaye, J.-Y., Barthès, L., Golé, P., Lavergnat J., and Vinson, J.P.: a dual beam spectropluviometer concept, Journal of
20  Hydrology, 328(1-2), 110-120, 2006.

Gargouri, E., Chebchoub, A.: Modélisation de la structure de dépendance hauteur-durée d'événements pluvieux par la copule de Gumbel. Hydrological Sciences–Journal–des Sciences Hydrologiques, 53(4), 802-817, 2010.

Grazioli, J., Tuia, D., and Berne, A.: Hydrometeor classification from polarimetric radar measurements: a clustering approach, Atmos. Meas. Tech., 8, 149-170, 2015.

25  Guyon, I. and Elisseeff, A.: An Introduction to Variable and Feature Selection  (Kernel Machines Section), 3, 1157--1182, 2003.Haile, A. T., Rientjes, T. H. M., Habib, E., Jetten, V., and Gebremichael, M.: Rain event properties at the source of the Blue Nile River, Hydrol. Earth Syst. Sci., 15, 1023-1034, 2011.

Iguchi, T., Kozu, T., Kwiatkowski, J., Meneghini, R., Awaka, J., and Okamoto, K.: Uncertainties in the rain profiling algorithm for the TRMM precipitation radar. J. Meteor. Soc. Japan, 87A, 1-30, 2009.

30  Holland, J. H.: Adaptation In Natural And Artificial Systems, University of Michigan Press, 1975.

Kohonen, T.: Self-organizing formation of topologically correct feature maps. Biological Cybernetics, 46, 59-69, 1982.

Kohonen, T. (2001). Self-Organizing Maps. Springer-Verlag, ISBN 3-540-67921-9, New York, Berlin, Heidelberg, 2001

35  Larsen, M. L. and Teves, J. B.: Identifying Individual Rain Events with a Dense Disdrometer Network, Advances in Meteorology, 2015, ID582782, 2015.

Lavergnat, J. and Golé, P.: A Stochastic Raindrop Time Distribution Model. Journal of Applied Meteorology, 37, 805-818, 1998.

Lavergnat, J. and  Golé, P.: A stochastic model of raindrop release: Application to the simulation of point rain observations,
40  Journal of Hydrology, 328(1), 8-19, 2006.

Liu, Y., Weisberg, R. H. and Mooers, C. N. K.: Performance evaluation of the self- organizing map for feature extraction, Journal of Geophysical Research, 111, C05018, doi:10.1029/2005JC003117, ISSN 0148-0227, 2006

Liu, Y. and Weisberg R.H.: A review of self-organizing map applications in meteorology and oceanography. In: Self-Organizing Maps-Applications and Novel Algorithm Design, 253-272, 2011.

Llasat, M.C.: An objective classification of rainfall events on the basis of their convective features. Application to rainfall intensity in the north east of Spain. International Journal of climatology, 21, 1385-1400, 2001.

Marzuki, M., Hashiguchi, H., Yamamoto, M. K., Mori, S., and Yamanaka, M. D.: Regional variability of raindrop size distribution over Indonesia, Ann. Geophys., 31, 1941-1948, 2013.

Molini, L., Parodi, A., Rebora, N., Craig, G. C.: Classifying severe rainfall events over Italy by hydrometeorological and dynamical criteria. Quarterly Journal of the Royal Meteorological Society, 137(654), 148-154, 2011.

de Montera, L., Barthes, L., and Mallet, C.: The effect of rain-no rain intermittency on the estimation of the Universal Multifractal model parameters, J. of Hydrometeorology, 10, pp. 493–506, 2009.

Moussa, R. and Bocquillon, C.: Caractérisation fractale d'une série chronologique d'intensité de pluie. Rencontres hydrologiques Franco-Romaines, 363-370, 1991.

Suh, S.-H., You, C.-H., and Lee, D.-I.: Climatological characteristics of raindrop size distributions in Busan, Republic of Korea, Hydrol. Earth Syst. Sci., 20, 193-207, 2016.

Tapiador, F. J., Checa, R., and de Castro, M.: An experiment to measure the spatial variability of rain drop size distribution using sixteen laser disdrometers. Geophysical Research Letters, 37(16), ID L16803, 2010.

Testud, J., S., Oury, P., Amayenc, and Black, R. A.: The concept of ''normalized'' distributions to describe raindrop spectra: A tool for cloud physics and cloud remote sensing, J. Appl. Meteor., 40, 1118–1140, 2001.

Ulaby, F. T., Moore, R. K., and Fung, A. K.: Microwave Remote Sensing: Fundamentals and Radiometry. Vol. I. Artech House, 321-327, 1981.

Uriarte, E. A. and Martín, F. D., Topology Preservation in SOM, World Academy of Science, Engineering and Technology. International Journal of Computer, Electrical, Automation, Control and Information Engineering 2, 9, 2008Verrier, S., Barthès L., Mallet C.: Theoretical and empirical scale dependency of Z-R relationships: Evidence, impacts, and correction, Journal of Geophysical Research: Atmospheres, 118 (14), 7435-7449, 2013.

Vesanto, J. and Alhoniemi, E.: Clustering of the self-organizing map. IEEE Transactions on Neural Networks, Vol. 11, 586–600, ISSN 1045-9227, 2000

# Tables

5

| | Observation period | Availability (%) | Number of rain events |
|---|---|---|---|
| Learning data set | 01/01/2013-12/31/2014 | 96.4% | 234 |
| Test data set | 04/16/2008-01/31/2012 | 60% | 311 |

**Table 1: Observation periods, availability of DBS observations, and numbers of rain events for the learning and test datasets**

| Number (#) | Feature name | symbol | Formula | Normalisation |
|---|---|---|---|---|
| 1 | Event duration | $D_e$ | $D_e = T_{end} - T_{begin} + 1$ [min] <br><br> With $T_{begin}$: Event start time and $T_{end}$: Event end time | 1 |
| 2 | Mean event rain rate | $R_m$ | $R_m = \frac{1}{D_e}\sum_{t=T_{begin}}^{t=T_{end}} RR_t$ [mm h⁻¹] | 2 |
| 3 | Intra-dry duration | $D_d$ | $D_d = \sum_{t=T_{begin}}^{t=T_{end}} I_t$ [min] <br><br> With $I_t = \begin{cases} 1 & if\ RR_t = 0\ [\text{mm } h^{-1}] \\ 0 & else \end{cases}$ | 0 |
| 4 | First quartile | $Q_1$ | The 25th percentile [mm h⁻¹] | 0 |
| 5 | Median | $Q_2$ | the 50th percentile [mm h⁻¹] | 0 |
| 6 | Third quartile | $Q_3$ | The 75th percentile [mm h⁻¹] | 2 |
| 7 | Previous IET | $IET_p$ | $IET_p = T_{begin}(current\ event) - T_{end}(previous\ event) + 1$ [min] | 0 |
| 8 | Mean rain rate over the rainy period | $R_{m,r}$ | $R_{m,r} = \frac{1}{(D_e - D_d)}\sum_{t=T_{begin}}^{t=T_{end}} RR_t$ [mm h⁻¹] | 3 |
| 9 | Event Rain rate std. | $\sigma_R$ | $\sigma_R = \sqrt{\frac{1}{D_e}\sum_{t=T_{begin}}^{t=T_{end}}(RR_t - R_m)^2}$ [mm h⁻¹] | 2 |
| 10 | Mode | $M_0$ | $M_0 = $ the most frequent RRₜ | 0 |
| 11 | Rain rate peak | $R_{max}$ | $R_{max} = \max(RR_t)$ | 2 |
| 12 | Dry Percentage in event | $D_{d\%e}$ | $D_{d\%e} = \frac{D_d}{D_e}$ | 5 |
| 13 | Rain event depth | $R_d$ | $R_d = R_m * D_e /60$ [mm] | 0 |
| 14 | Absolute rain rate variation of order c | $P_{c1}$ | $P_{c_i} = \sum_{t=T_{begin}}^{t=T_{end}-1} |RR_{t+1} - RR_t|^{c_i}$ <br><br> For $c_i = 0.5, 1, 2$ | 6 |
| 15 |  | $P_{c2}$ |  | 3 |
| 16 |  | $P_{c3}$ |  | 2 |
| 17 | Normalized Absolute rain rate variation of order cᵢ | $P_{C_{N1}}$ | $P_{C_{Ni}} = \frac{P_{c_i}}{D_e}$ <br><br> For $i = 1..3$ | 3 |
| 18 |  | $P_{C_{N2}}$ |  | 2 |
| 19 |  | $P_{C_{N3}}$ |  | 0 |
| 20 | Absolute rain rate variation of order C and threshold S | $P_{S,C}$ | $P_{S,C} = \sum_{t=T_{begin}}^{t=T_{end}-1} |\max[(RR_{t+1} - S), 0] - \max[(RR_t - S), 0]|^C$ <br><br> With $s = 0.3$ and $c = 2$ | 6 |
| 21 | $\beta_L$ parameter | $\beta_{L1}$ | $\beta_{L_i} = \frac{\sum_{i=T_{begin}}^{T_{end}} RR_t\ \theta(RR_t - L_i)}{\sum_{i=T_{begin}}^{T_{end}} RR_t}$ <br><br> For $L_i = 0.3, 1, 3\ mm\ h^{-1}$ <br><br> With $\theta(RR_t - L_i)$ is the Heaviside function defined as <br> $\theta(RR_t - L_i) = 1\ if\ RR_t \geq L_i$ <br> $\theta(RR_t - L_i) = 0\ if\ RR_t < L_i$ | 5 |
|  |  | $\beta_{L2}$ |  | 0 |
| 22 |  |  |  |  |
| 23 |  | $\beta_{L3}$ |  | 0 |

**Tablec2: The 23 variables identified in the literature, used for the characterization of rain events**

| Transformation number | Transformation name | Formula f(x) | Note & remark |
|---|---|---|---|
| 0 | Standardisation | $\dfrac{x - mean(x)}{std(x)}$ | - |
| 1 | Power | $x^n$ | $n = 0.05$ |
| 2 | Boxcox | $\dfrac{(x^\gamma - 1)}{\gamma}$ | $\gamma = -0.1$ |
| 3 | | | $\gamma = -0.2$ |
| 4 | | | $\gamma = -0.3$ |
| 5 | Arc-sin of square | $\arcsin \sqrt{x}$ | Data are between 0 and 1 |
| 6 | decimal Logarithm | $Log(x + c)$ | $c = 0.1$ |

**Table 3: Transformations used to normalize the variables listed in Table 2.**

| Variables | $D_e$ | $R_m$ | $D_d$ | $Q_1$ | $Q_2$ | $Q_3$ | $IET_p$ | $R_{m,r}$ | $\sigma_R$ | $M_0$ | $R_{max}$ | $D_{d\%e}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $R^2$ learninglearning data set | 0.96 | 0.91 | 0.58 | 0.57 | 0.48 | 0.77 | 0.31 | 0.93 | 0.97 | 0.50 | 0.96 | 0.52 |
| $R^2$ Test data set | 0.93 | 0.84 | 0.55 | 0.57 | 0.50 | 0.74 | 0.26 | 0.84 | 0.86 | 0.54 | 0.82 | 0.50 |

| Variables | $R_d$ | $P_{c1}$ | $P_{c2}$ | $P_{c3}$ | $P_{cN1}$ | $P_{cN2}$ | $P_{cN3}$ | $P_{S,C}$ | $\beta_{L1}$ | $\beta_{L2}$ | $\beta_{L3}$ | - |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $R^2$ learning data set | 0.97 | 0.97 | 0.93 | 0.94 | 0.91 | 0.95 | 0.78 | 0.94 | 0.70 | 0.89 | 0.96 | |
| $R^2$ Test data set | 0.94 | 0.99 | 0.91 | 0.83 | 0.83 | 0.85 | 0.71 | 0.82 | 0.61 | 0.76 | 0.89 | - |

**Table 4: Coefficient of determination obtained on the learning and test datasets. The values with a dark grey background correspond to the 5 selected variables**
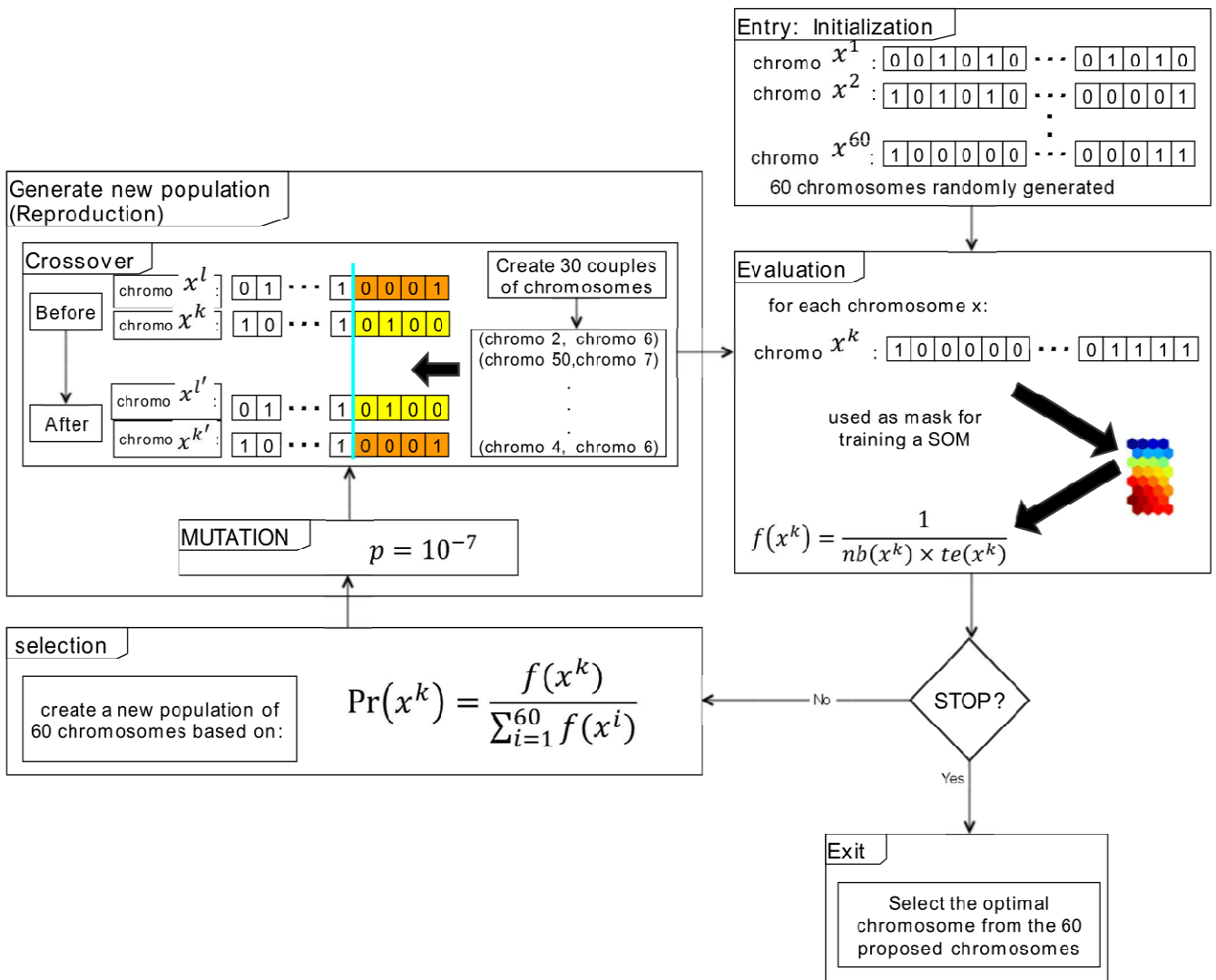
| | Stratiform events | | Convective events | | |
|---|---|---|---|---|---|
| | **Subclass 1** | **Subclass 2** | **Subclass 3** | **Subclass 4** | **Subclass 5** |
| Variables | Mean | Mean | mean | mean | mean |
| $D_e(min)$ | 321 | 149 | 464 | 75 | 49 |
| $\sigma_R$ | 0.36 | 2.01 | 3.62 | 11.7 | 9.64 |
| $R_{max}(mm\ h^{-1})$ | 2.08 | 10 | 22 | 52.7 | 36.06 |
| $R_d(mm)$ | 1.99 | 2.62 | 11.24 | 6.9 | 2.72 |
| $P_{c1}$ | 75.7 | 64.5 | 193 | 78.2 | 40.94 |
| $R_m(mm\ h^{-1})$ | 0.37 | 1.48 | 2.35 | 7.85 | 7.11 |
| $D_d(min)$ | 80 | 31 | 75 | 11 | 1 |
| $\beta_{L3}$ | 0.01 | 0.42 | 0.48 | 0.89 | 0.86 |

**Table 5: Summary of the rain event subclasses computed with the learning dataset.**
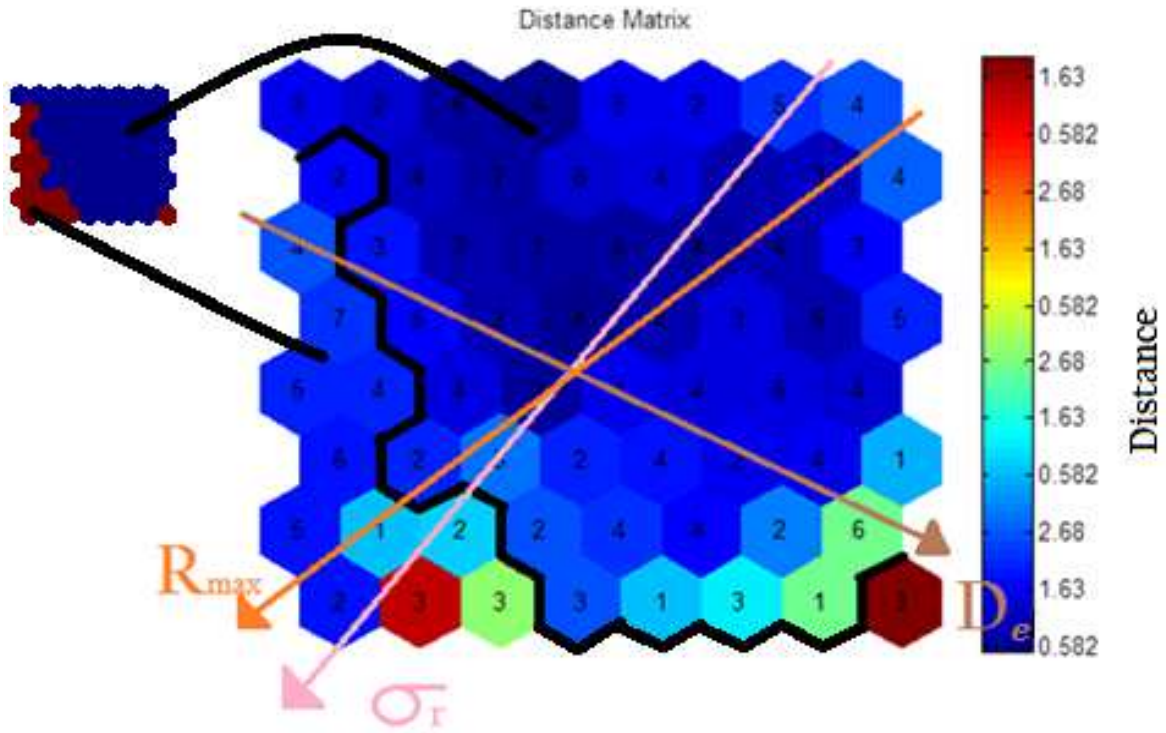
# Figures





**Figure 1: PCA on the learning data set based on the 23 variables described in Table 3. Left: correlation circle on axes 1 & 2. Right: correlation circle on axes 1 & 3. All of the variables are normalised according to the last column of Table 2.**
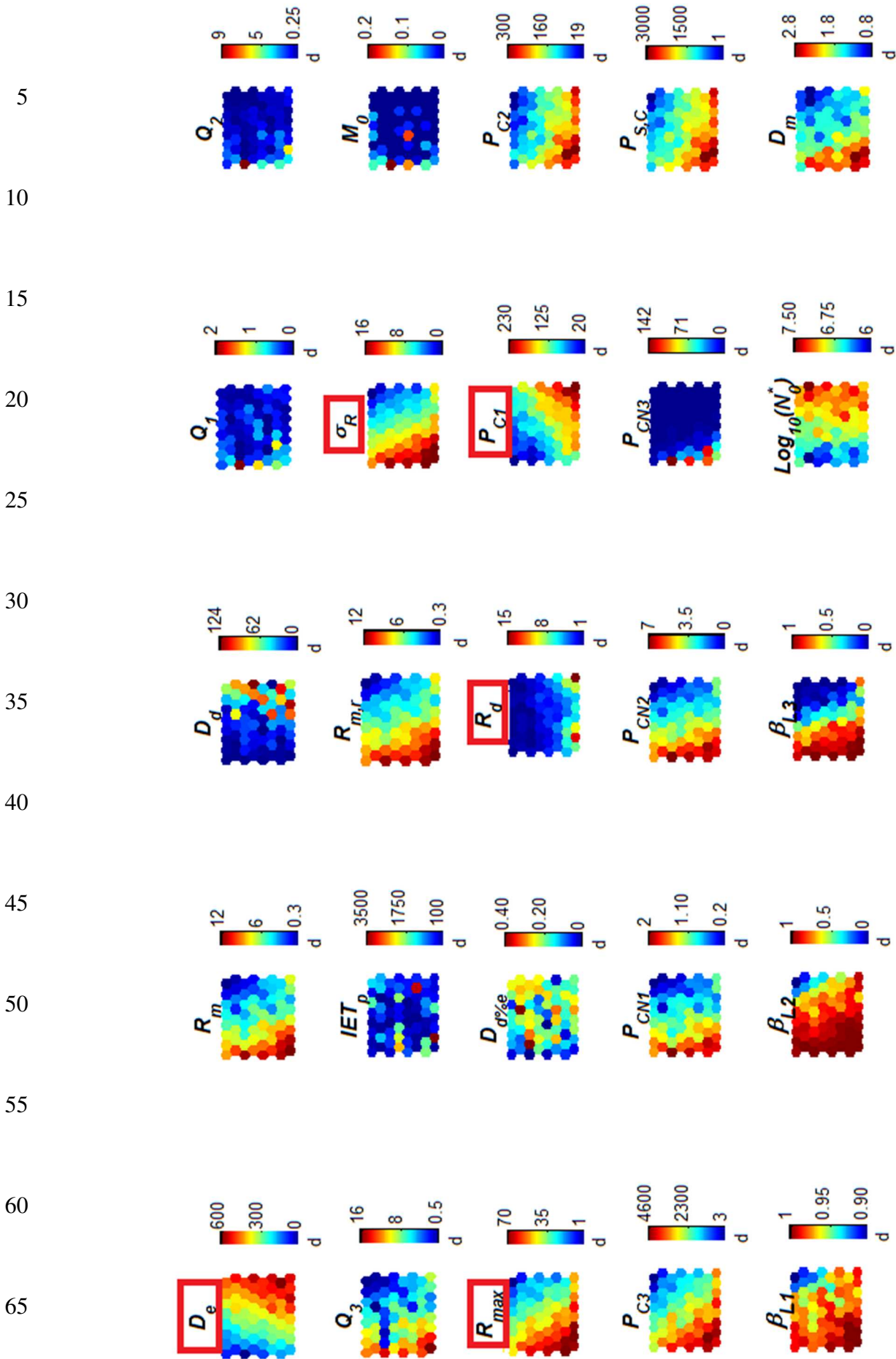
**Figure 2: Diagram for the selection of variables based on a Genetic Algorithm associated with Kohonen Maps**
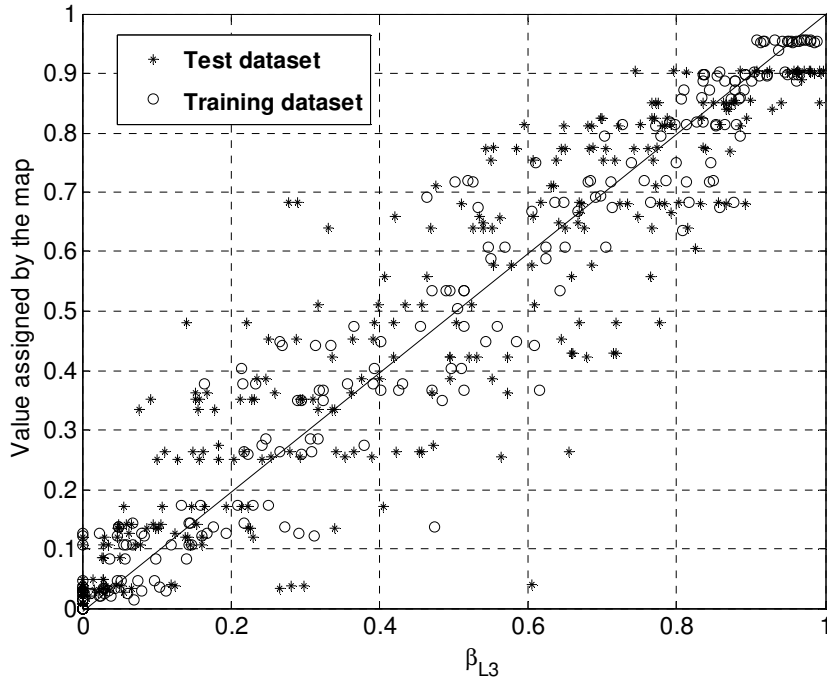
Figure 3: Distance matrix for the $M(x^{Best})$ map: The colour of each neuron represents the average distance between itself and its neighbouring neurons. The value inside each neuron indicates the number of rain events that it has captured inside the learning dataset. The black line separates the neurons into 2 classes, using the Hierarchical Ascendant Classification (see section 4.1). The arrows represent the gradients of the variables $R_{max}$, $\sigma_R$ and $D_e$

**Figure 4: Projection of the $M(x^{Best})$ map according to the 23 variables. The red-framed variables are those selected by the GA algorithm. The last two variables $D_m$ and $N_0^*$ are defined in section 5.**

**Figure 5: The variable $\beta_{L3}$ versus its corresponding value, given by the best matching unit: from the learning dataset (circles), and the test dataset (stars). The solid line corresponds to the first diagonal.**
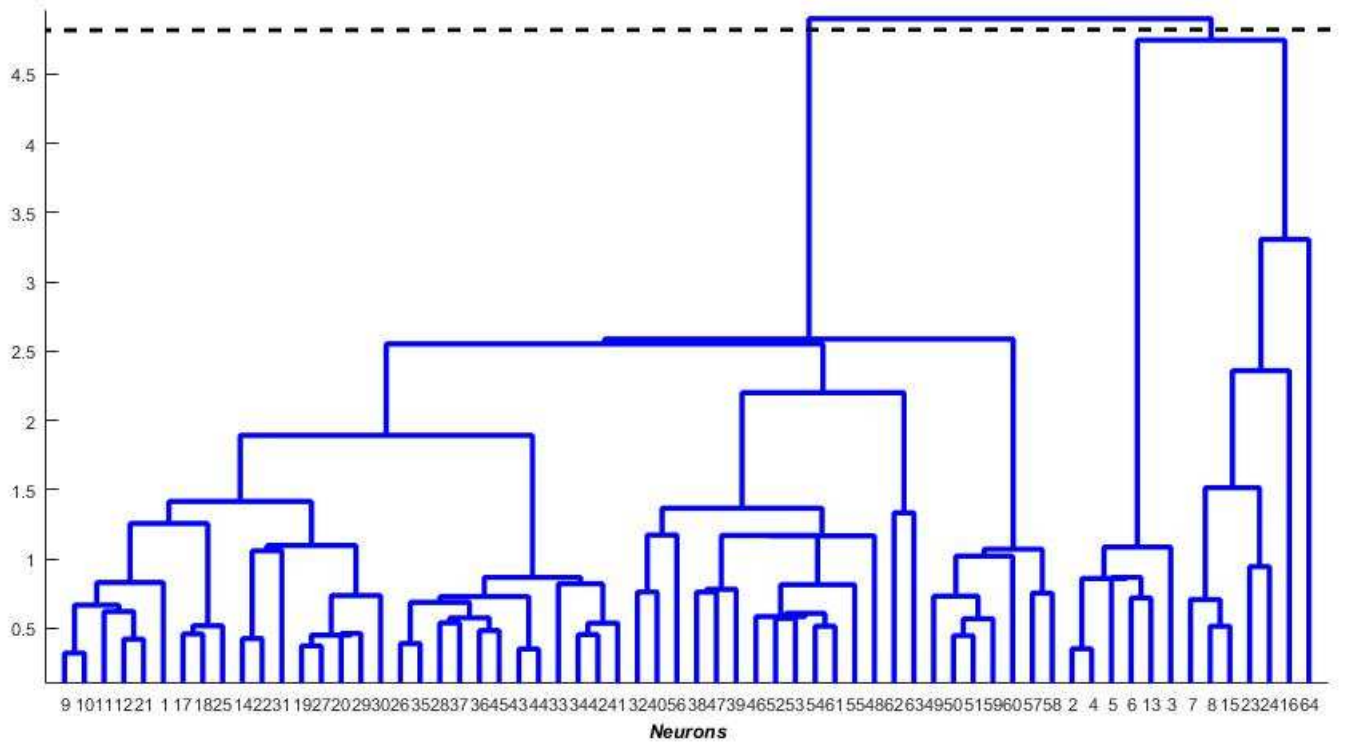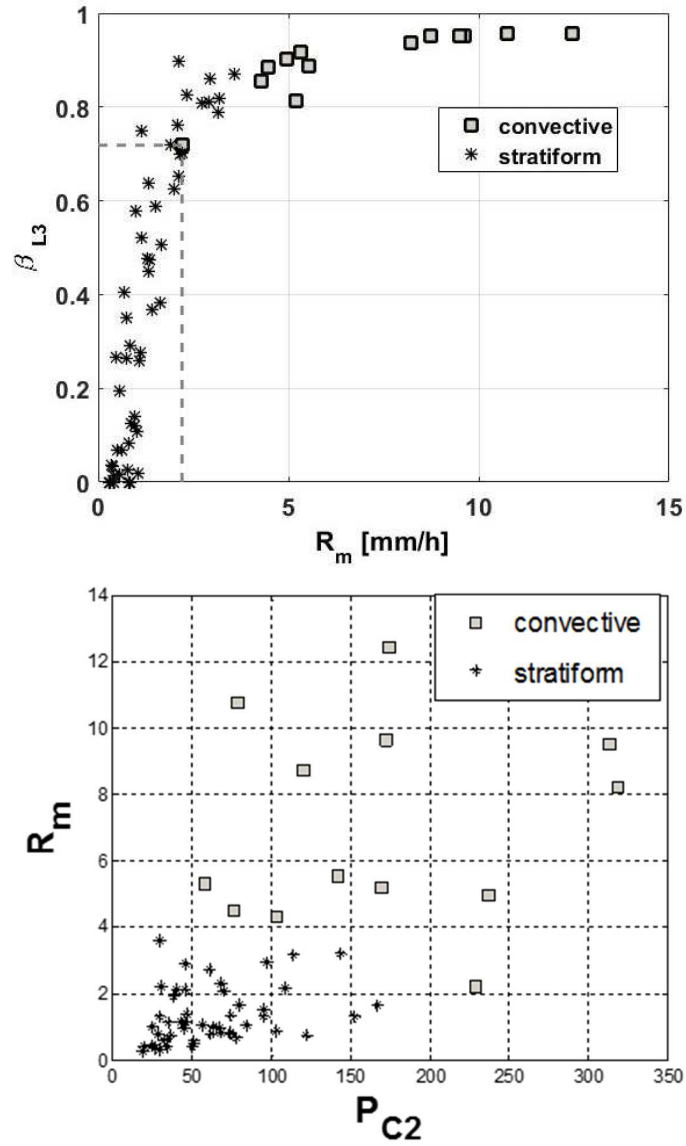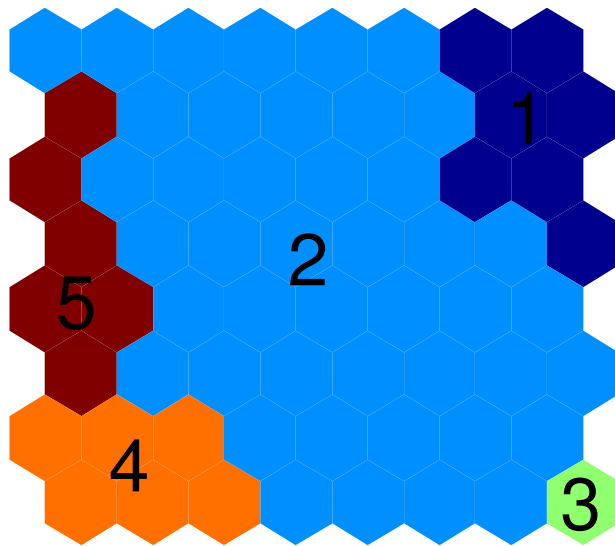
**Figure 6  Dendrogram obtained from the Hierarchical Cluster Analysis of the 64 neurons in the SOM. The horizontal dashed line represents the threshold between the two classes.**
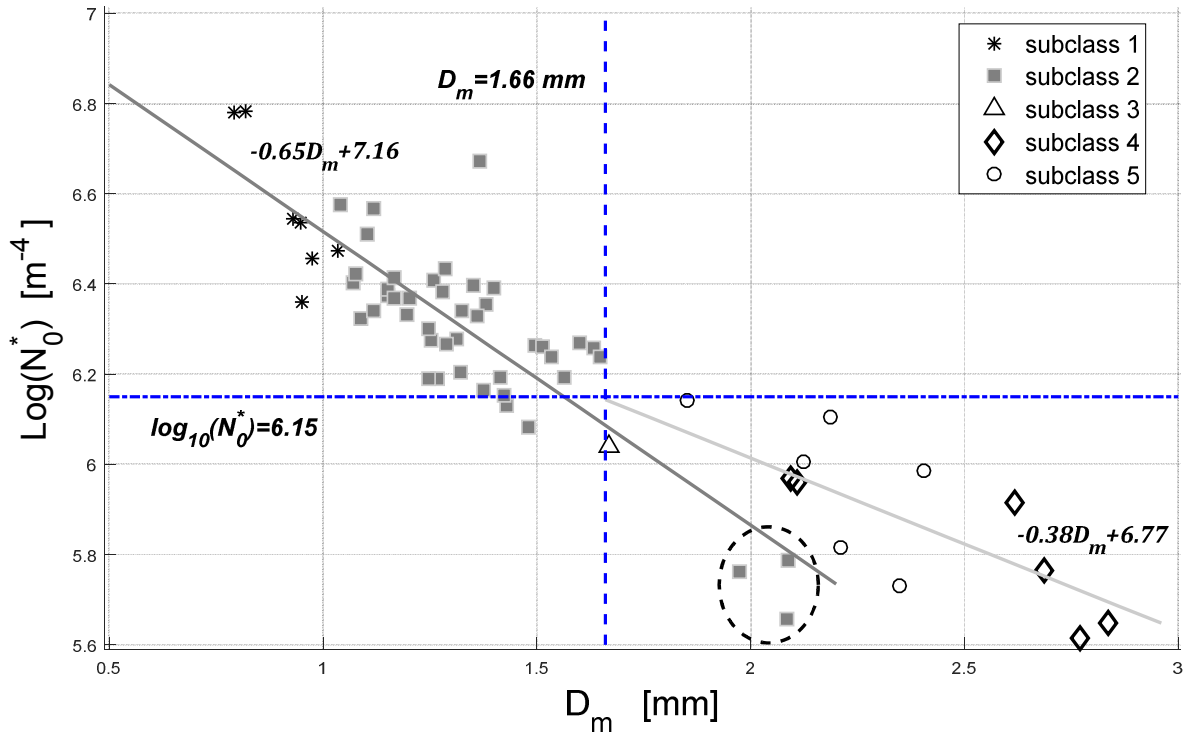
5

5    Figure 7: Representation of the neurons in the Rm, $\beta_{L3}$ and Rm, $P_{C2}$ sub-spaces. The stars represent neurons from group 1 (stratiform), and the squares correspond to neurons from group 2 (convective). Dashed lines indicate the neuron #64.

29

**Figure 8: Hierarchical Clustering of the map into five subclasses. The colours represent the subclass numbers: Subclass 1: Dark blue, Subclass 2 : blue, Subclass 3 : Green, Subclass 4 : Orange, Subclass 5 : Red**

5

*Figure 9: Microphysical variable $N_0^*$ versus $D_m$ for the five rainy event subclasses. The three neurons corresponding to mixed events are circled. The dashed lines correspond to borders $D_m > 1.66$ and $Log(N_0^*)>6.15$*