



# A Multilevel Clustering Method for Risky Areas in the Context of Avalanche Danger Management

Fanny Pagnier, Frédéric Pourraz, Didier Coquin, Gilles Mauris, Hervé Verjus

## ► To cite this version:

Fanny Pagnier, Frédéric Pourraz, Didier Coquin, Gilles Mauris, Hervé Verjus. A Multilevel Clustering Method for Risky Areas in the Context of Avalanche Danger Management. Information Processing and Management of Uncertainty in Knowledge-Based Systems, 1602, Springer International Publishing, pp.54-68, 2022, Communications in Computer and Information Science, 10.1007/978-3-031-08974-9\_5 . hal-03728237

**HAL Id: hal-03728237**

**<https://hal.science/hal-03728237>**

Submitted on 20 Jul 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A multilevel clustering method for risky areas in the context of avalanche danger management

Fanny Pagnier<sup>1</sup>, Frédéric Pourraz<sup>1</sup>, Didier Coquin<sup>1</sup>, Hervé Verjus<sup>1</sup>, and Gilles Mauris<sup>1</sup>

LISTIC - Université Savoie Mont Blanc, 5 chemin de Bellevue, Annecy-le-Vieux,  
74940 Annecy, France  
`{surname.name}@univ-smb.fr`

**Abstract.** In the context of avalanche risk management, we study the spatial variability of rainfall conditions, which is one of the main parameters that induce natural avalanches. This paper focuses on the geographical variability of the snow overload due to recent precipitations. We propose a generic approach applicable at a larger scale and, for this reason, without relying on any expert knowledge. Our proposal is a multilevel clustering process based on classical methods processed in sequence to take advantage of each one. As a result, the multilevel clustering process outputs four main detected weather trends that affect the French Alps. The developed process is generic enough to be used in other areas. Our work is intended to positively impact and improve the current and future decision support methods and tools for mountain practitioners.

**Keywords:** Clustering · Unsupervised methods · Decision making .

## 1 Introduction

In avalanche risk management, avalanche observations are a particularly effective indicator of the current danger level and a direct evidence of snow instability [7]. Avalanche observations are indeed one of the input factors of several decision support methods. Until now, this factor remains considered as a warning sign and is never quantified or measured in these methods. The reader can refer to the most recent survey considering the factors and methods used by experts [6]. For example, the Obvious Clues Method [10] remains vague by solely considering “*Avalanches in the area in the last 48h*” but not formalizing the boundaries of the area under consideration. However, to take correct decisions, it is beneficial to estimate where similar conditions as the one leading to the observed avalanche are likely to be encountered. This paper proposes a first approach to quantify the size of the area to take into consideration. Although several parameters can be taken into account to consider this notion of similar conditions, the paper focuses on the amount of accumulated snow during the last 24 hours.

In this general context, the main objective is to identify which French Alps areas are likely to receive similar amounts of fresh snow, according to different weather trends. It is equivalent to studying the variability in rainfall conditions at lower altitudes, where the density of automatic measuring stations is higher. Such studies and proposed work in this paper are intended to positively impact and to improve the mountain practitioners’ decision support methods and

tools. Scientific developments in avalanche science adopt a multidisciplinary approach covering field observations and experiments as well as mathematical and physical analysis modeling [9]. Most researchers are interested in the internal snowpack variability at various scales [16] but experts consider the snowpack’s overload as one of the main criteria causing natural avalanches. For this reason, this paper focuses on the geographical variability of the snow overload due to recent precipitations. This work is part of a French-Swiss Interreg project (see Acknowledgements) and is currently being evaluated at the French Alps scale. We aim to propose a generic approach applicable at a larger scale and, for this reason, without relying on any expert knowledge. Several works [8][1][5] have already addressed the study of precipitations in different countries and regions but there is no general method that could be applied on a wide scale. While two locations may encounter the same rainfall conditions under a given weather trend, they may behave differently under a different one. Moreover, some areas may be affected by several weather trends. Thus, when observing a new natural avalanche, it is necessary to estimate the current weather trend of this specific day to know which area makes sense for the given observation.

The specific objective of this study is thus to classify days that are similar in terms of the location of the main rainfall totals. This way, we bring out the major weather trends affecting a given mountain area and we can determine the influencing zone of an avalanche observation.

We work on a dataset based on 12 years of rainfalls measurements collected during the winter season. These measurements are recorded hourly from 90 stations spread on the French Alps. Data comes from EDF-DTG’s<sup>1</sup> automatic measuring stations.

Early, we understood that applying clustering methods on the whole dataset does not give the expected result. It classifies days according to the total amount of new rainfalls instead of depending on the location of the main rainfall totals. Figure 1 illustrates this assessment, showing three clusters with the same rainfalls location but various intensities. Then, to address this problem, we decided to work on subsets, which considerably improve the results: when reducing the size of the dataset (i.e., the number of individuals), we noticed that the clustering better captured the location of the precipitations rather than focusing on the total amount of precipitations. Applying clustering methods on 12 separated data subsets implies merging the results obtained on each separated subset into a final global one. That is why we then develop a two-level clustering process.

The paper is structured as follows: Section 2 presents step-by-step the developed process and the methods used. Then, Section 3 shows the obtained results and validation cases. We discuss in Section 4 the developed process and give some perspectives relative to this work. Finally, Section 5 conclude the paper.

## 2 Multilevel clustering method

Machine learning approaches can be classified into supervised learning, unsupervised learning, and reinforcement learning. Since we do not have any labeled data

---

<sup>1</sup> *Électricité de France*, the French electricity production and supply company

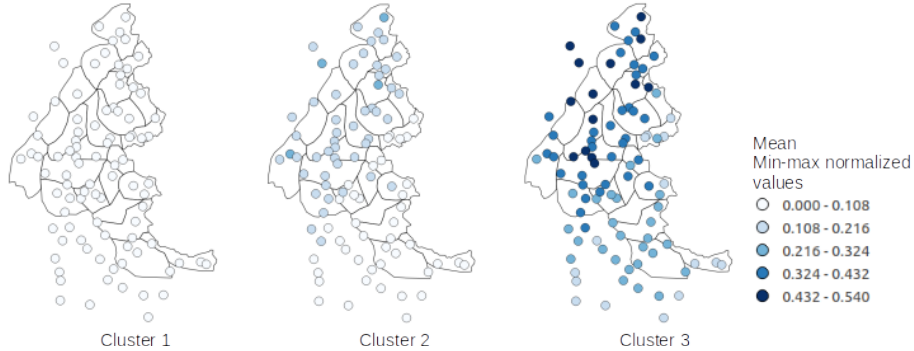


Fig. 1: Clustering (PCA  $\rightarrow$  HC  $\rightarrow$  K-means) on the entire dataset (i.e., 12 winter seasons). For details about the use of the sequence PCA  $\rightarrow$  HC before K-means, see Section 2.1.

and we can not do trial-and-error experiments in the context of risk management, we focus on unsupervised methods. We present a classification approach that 1) meets our previously stated objectives, 2) is an automatic process that does not require any parameter refinement nor specification, and 3) is a transferable process working end-to-end without expert assistance. For this purpose, we first use classical methods of statistical analysis and machine learning: Principal Component Analysis (PCA) [15], Hierarchical Classification (HC) [14], and K-means [4][19] successively, and to refine the result, we use Affinity Propagation (AP) [3].

The whole multilevel clustering process is presented in Figure 2. In addition, Figure 4, which illustrates the whole process once applied to our data (see Section 3), gives also a general overview of the multilevel clustering process.

For our problem, we consider a dataset structured as follows: individuals are days, and variables are rainfall totals recorded over 24 hours in a given meteorological station and for a given day. We thus manage both temporal and spatial components. The temporal component, first, makes it possible to split the initial dataset into 12 data subsets (to work on each winter season separately). Then, as mentioned here, the temporal component permits us to identify the studied individuals : we classify days. For each individual, the variables correspond to measures of rainfall totals received in 90 stations spread over the French Alps, which induces an indirect spatial component as every measuring station corresponds to a specific location on the territory. But, the spatial component relative to the geographical location of the measuring stations is never taken into account (we do not consider information as latitude or longitude). However, this emerges due to PCA whose main components reflect the geographical location of the stations (see Section 2.1 and [13] for more details). In addition, the spatial component is visible when visualizing the results on maps. We will see that results are consistent with the spatial location of the measuring stations. To summarize, the clustering process aims to group individuals (i.e., days) that are similar in

terms of the location of the main rainfall totals (i.e., days receiving the main rainfall totals on the same measuring stations). At the end of the first-level clustering process, each cluster (i.e., group of days) corresponds to a meteorological trend (all days received similar amounts of precipitations at the same location). That finally permits the identification of geographical extents that are impacted in the same way according to a given trend.

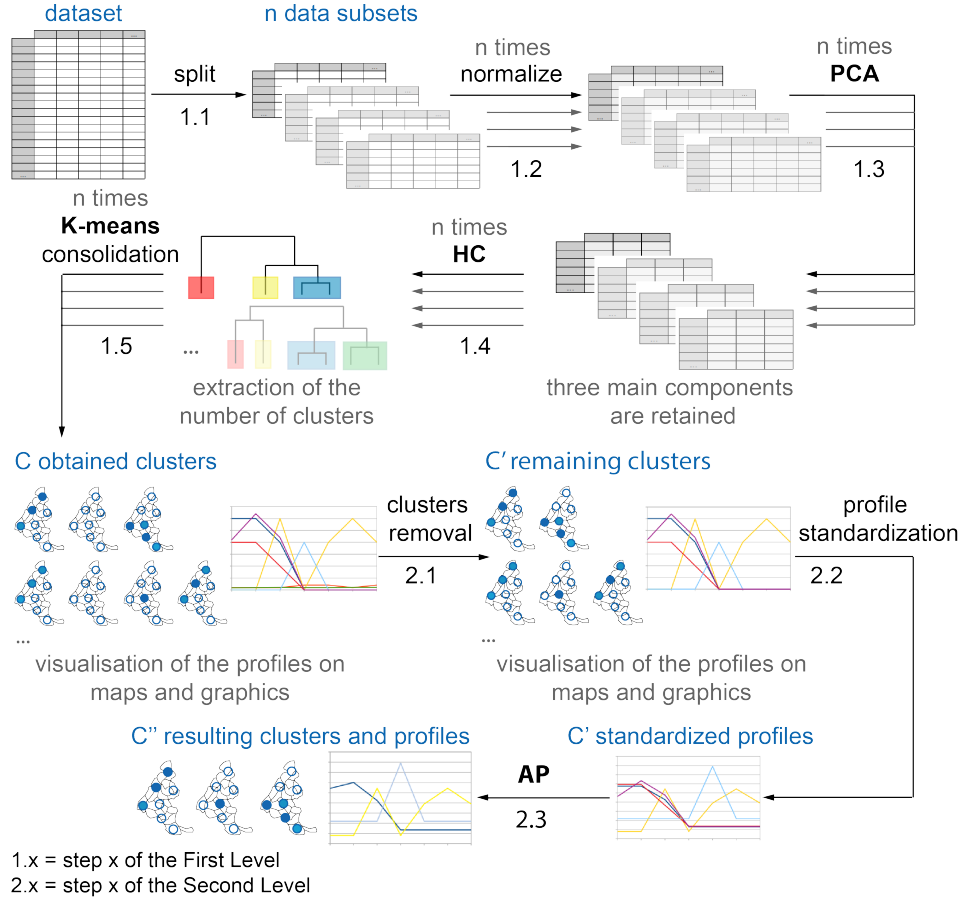


Fig. 2: Schematic representation of the multilevel clustering method

## 2.1 First Level (PCA → HC → K-means)

The first level contains 5 successive steps. First, for the reason mentioned in the introduction, the initial dataset is split into subsets (see step 1.1 on Figure 2). Here, we decide to process each winter season separately. Thus we obtain  $n$  data subsets. Each data subset contains  $x$  individuals and  $y$  variables ( $x$  and  $y$  can be slightly different depending on the processed year).

Data subsets have a feature that requires pre-processing (see step 1.2 on Figure 2): the variables have different orders of magnitude whereas all of them need

to have the same importance for the following treatments. As pre-processing, we normalize the range of the variables to values between 0 and 1 and give them all the same importance for the subsequent treatments. All values are divided by the maximum value recorded on the station over the 12 seasons (i.e.. the whole available data). That corresponds to a min-max normalization. In our context, for all variables:  $min = 0$ . So, the equation becomes:

$$value/max \quad (1)$$

All variables will translate the idea that the recorded rainfalls over the last 24 hours are small or high according to the maximum rainfall total the station have ever received in the last 12 years. This pre-processing allows working on the importance of the received rainfalls relative to each station's features.

After pre-processing, the main part of the first level consists of three methods used in sequence to obtain a clustering result [12][17]. These three methods are: PCA [15], HC [14] and K-means [4][19]. This sequence corresponds to steps 1.3, 1.4 and 1.5 on Figure 2.

The goal of **PCA** is to transform the  $y$  variables into a reduced number of principal components. Each component corresponds to a linear combination of the initial variables, and only the firsts are truly informative as they contain most of the variability. Then, to work on a reduced number of components reduces the problem dimension, simplifies results in interpretation after classification, and removes a part of the noise contained in the data subset due to the nature of the studied phenomenon. At the end of the PCA stop, we select only three principal components as they capture most of the initial variability and are well explainable. The first one corresponds to the opposition between dry and wet days, whereas the second and third components correspond to North / South and East / West oppositions. More details are given in [13].

After the PCA, HC and K-means are run in sequence to take advantage of both methods. **HC** creates clusters by aggregating elements two by two. During the successive iterations, the method creates from  $x - 1$  clusters to a sole cluster (which contains all the  $x$  individuals). We use the euclidean distance and Wards criterion [18], which minimizes the loss of between-clusters inertia when aggregating two clusters. Thanks to the **HCPC** function from the **FactoMineR** package of R, the number of clusters the most appropriate for our data subset is automatically obtained. The partition is the one with the biggest relative loss of inertia. We call this number of clusters  $c_i$  (where  $i$  represents each data subset, i.e.,  $i \in [1, n]$ ). This method has, however, a drawback: when it misclassifies an individual, it remains misclassified until the end. Using the **K-means** algorithm improves the classification. It offers each individual the possibility to move from one cluster to another during successive iterations. The K-means algorithm has to be initialized, which is possible thanks to the result given by the HC and avoids a random initialization. As an objective was to develop a generic process suitable for other data (for example, in the Swiss Alps) without any expert knowledge, we looked for a method suggesting a first estimation of the  $k$  value. This way, K-means is not set randomly but corresponds to a consolidation of the HC's result.

This sequence (PCA  $\rightarrow$  HC  $\rightarrow$  K-means) is applied to each data subset. At the end of the first level clustering, it gives in total  $C$  clusters, where  $C = \sum_{i=1}^n c_i$ . The number of clusters ( $c_i$ ), automatically estimated by the process, may differ on each data subset. Hence, here are some drawbacks:

- The process does not always detect every weather trend, even if they are effectively present in the data subset;
- The process sometimes returns several clusters corresponding to the same weather trend (i.e., same location of the precipitations) but with different intensities.

These drawbacks explain why we use a multilevel clustering process:

- Repeating the first level clustering on several data subsets (so on several winter seasons) increases the possible detection of all trends;
- A second level of clustering is required to group similar clusters (i.e., clusters that correspond to the same area impacted by the rainfalls but with various intensity), whether they are obtained on one data subset or among the different winter seasons.

Of course, at the end of the whole process, the aim is to detect all the main trends, and to obtain only one final cluster per trend, i.e., to classify together all the clusters that have similar rainfalls locations but different intensities. To explain that, we introduce the notion of *rainfalls profile*.

The rainfalls profile of a day is the graphical representation of the received rainfalls on each of the 90 measuring stations (i.e., on the 90 variables). This profile can be represented on a graphic or a map (according to the location of the stations), and permits the identification of the location the most affected by the rainfalls. The profile can represent the raw data or the normalized ones (given after equation 1). Figure 3 shows, for the 08th January 2018, the normalized profile on a graphic (Figure 3-a, values between 0 and 1) and the raw profile on a map (Figure 3-b, values corresponding to rainfalls in mm). The most affected stations (and by extension areas) correspond to peaks on the graphic or to dark points on the map. The profile of a cluster (which is a group of days) corresponds to the mean of the days' profiles.

Profiles that have similar shape (on the graphic) correspond to the same trend as the same area is impacted. For a given trend, profiles may vary in values (that corresponds to more or less intense rainfalls). As the need is to classify together all the clusters that have similar rainfalls locations even if they have different intensities, the process needs to group profiles that are similar in shape even if they are different in values. Indeed, a given weather trend impacts the same areas, whatever the rainfall total.

In this context, for the second level clustering described below, we now consider the clusters profiles as new individuals. As first level gives  $C$  clusters,  $C$  different profiles are considered. Clusters' centroids are not taken as new individuals for the second level because the values on the 3 components given by PCA are in a specific referential given by the calculation of the PCA on each data subset independently. Coming back to all the initial variables makes it possible to have a common referential for all clusters.

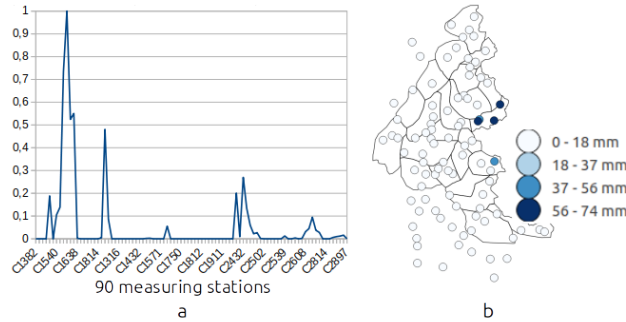


Fig. 3: Illustration of the notion of *rainfalls profile* (08th January 2018)

## 2.2 Second Level (AP)

At this stage, the clusters composed of days with no or few precipitations can be removed. Indeed, as we are working in the context of avalanche risk, there is no interest for these clusters as they do not correspond to an important enough overload to impact the snowpack's stability and generate avalanche activity. We decide to remove all the clusters for which the profile is, on each variable, lower than 0.2 (see step 2.1 on Figure 2). Note that, when removing days with no or few precipitations (i.e., values lower than 0.2 on each variable) before starting the multilevel clustering method, the first level process outputs anyway clusters corresponding to few rainfalls (i.e., profile lower than 0.2 on each variable, because values upper than 0.2, which induce to keep these days, are sparse, and not on the same stations). Thus, the sole deletion of clusters, at step 2.1, allows a sole data deletion. In addition, it makes it possible to keep the most available information of the initial data for the PCA step in order to better highlight the spatial components described above. As variables have been normalized (see above step 1.2), 0.2 means that precipitations were less than 20% of the maximum rainfall total observed in 24 hours on the station over the 12 seasons. This threshold is arbitrary, but some experimentations with 0.1 and 0.3 thresholds were carried out and variations had no impact on the final result. At this stage,  $p$  profiles are removed. It remains then  $C'$  profiles considered as individuals for the second level, where  $C' = C - p$ .

Then, as the need is to group profiles similar in shape but different in intensity, we have to standardize (i.e., center and reduce) them before starting a new clustering method. This way, profiles similar in shape become also more similar in values (see step 2.2 on Figure 2). For each profile, values are replaced by:

$$(value - \mu)/\sigma, \quad (2)$$

where  $\mu$  is the mean and  $\sigma$  the standard deviation of the 90 values that composed the profile.

The  $C'$  standardized profiles can then be classified. For this clustering level (see step 2.3 on Figure 2), we decide to use the **AP** method [3] as it does not



require setting the number of clusters nor choosing an initial set of points. This algorithm considers each individual as a potential exemplar and uses a similarity matrix between all individuals two by two. This method fits our objectives as it aims to classify together similar profiles.

As presented in the literature from other fields [11], the AP’s unsupervised version (which consists in using the preference value by default, i.e., the median value of the similarity matrix) usually gives too big of a number of clusters [3]. It was our case when we tried to use AP directly after PCA at the first level clustering (instead of HC and K-means). But here, as we are working on pre-processed data and mean individuals obtained through the first level clustering, the default preference value (*pref*) gives an acceptable number of clusters in the output. The ratio between *pref* and the maximum value of the similarity matrix is approximately 40%. This ratio was approximately 1% when using AP at the first level clustering. This difference explains the former observation. The higher *pref* is, the lower is the number of final clusters. So, at the second level clustering, the value of *pref* set by default is adapted to give a satisfying number of clusters in the output.

Thus, thanks to AP at this second level, the process finally gives a single cluster per weather trend. It correctly returns all the detected weather trends spotted during the previous level. The process finally returns  $C''$  final clusters. Each resulting cluster corresponds to a group of profiles similar in shape (i.e., profiles that correspond to similar rainfalls locations). Then, at this stage, we can visualize on maps the resulting profiles (which are means of standardized profiles processed at this second level). This gives a first idea of the rough areas mainly affected by rainfalls according to each detected weather trend.

### 3 Results and validation

#### 3.1 Multilevel clustering applied to our data: results

This section gives the result of the whole process  $\text{PCA} \rightarrow \text{HC} \rightarrow \text{K-means} + \text{AP}$  applied to our dataset (Figure 4).

Our dataset contains 12 winter seasons (i.e., days from 1st of December to 31th of March) from 2009-2010 to 2020-2021. That corresponds to a sum of 1455 individuals. Once split, there are 12 data subsets (i.e.,  $n = 12$ ), each corresponding to a winter season. According to the winter season, data subsets contain 121 or 122 individuals and 90 variables.

We pre-process each data subset and apply the first level clustering. When keeping the three first principal components of the PCA, we keep between 84.5% and 90.5% of the initial variability of the data subsets. HC suggests generating between three and seven clusters (i.e.,  $c_i \in [3; 7]$ ) per the data subset. We finally apply K-means to consolidate results. This way, on the 12 data subsets, we obtain 48 clusters in total (i.e.,  $C = 48$ ). We then calculate the 48 corresponding profiles.

Before the second level clustering, we pre-process the data by removing the profiles corresponding to a few rainfalls (with a 0.2 threshold), that is to say

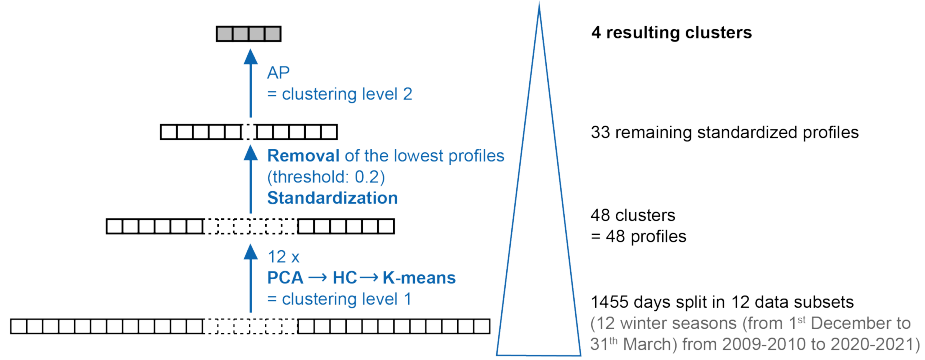


Fig. 4: Multilevel clustering applied to our data

here 15 profiles (i.e.,  $p = 15$ ). We then calculate the reduced profiles on the 33 remaining ones (i.e.,  $C' = 33$ ). Then, we apply the second level clustering, i.e., AP. The process gives at the end four clusters (i.e.,  $C'' = 4$ ) as result. Figure 5 shows their profiles on maps. The four resulting clusters correspond to four fictitious days, which are typical of the four main weather trends that the process detects (Figure 5). These trends are:

1. Most impacted area located in the southern part of the French Alps that corresponds to a flux mostly coming from South or South-West directions;
2. Most impacted area located in the northern part of the French Alps that corresponds to a flux coming from North-West;
3. Most impacted area located in the eastern side of the French Alps that corresponds to a flux coming back from the East;
4. Most impacted area located in the western part of the French Alps that corresponds to a flux coming mostly from the West, affecting, the pre-alpine massifs.

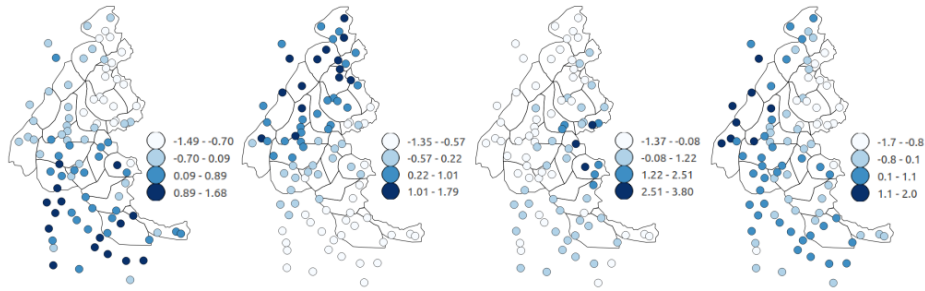


Fig. 5: Visualization of the reduced mean profiles of the resulting clusters

We presently know that at least four different meteorological trends may be detected in the 12 winter seasons dataset. Let us see which result gives K-means algorithm directly leads on the three main components extracting by PCA on the

whole dataset (i.e., 12 winter seasons), with random initialization and  $k = 4$ . In this case, we detect two different trends (see Figure 6, clusters 1 and 4 correspond to two different trends, with the same range of intensity) and one of these trends appears with three different intensities (see clusters 2, 3, and 4 on Figure 6). Thus, applying K-means on the entire dataset does not give an optimum result (as solely two out of the four possible trends are detected, and as a sole trend is given several times): even by imposing  $k = 4$ , it does not output the four possible trends. It emphasizes that the multilevel clustering process (working on a split dataset in the first level and merging the results in the second one) produces better outputs.

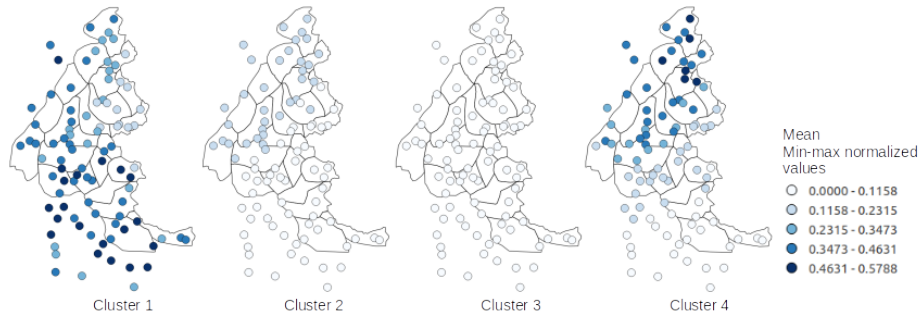


Fig. 6: Clustering (PCA  $\rightarrow$  K-means) on the entire dataset (i.e., 12 winter seasons), with  $k=4$

### 3.2 Two validation cases

The process' strength is its global result (obtained through two levels clustering and the calculation of means to generate the new individuals processed in the second level). That means that each initial individual weights less than the collective does: even if some individuals are misclassified (because the first level does not detect every weather trend and does not strictly classify each individual well), they do not distort the final global result.

The two following examples illustrate this outcome. Experts identify the 08th and 09th January 2018 as two days affected by a flux coming from the east. Under this weather trend, precipitations are located mostly on the eastern part of the French Alps (Figure 7 illustrates this assumption).

The process misclassified these two individuals. The location of the rainfalls of the 09th of January 2018 (Figure 7) differs from the profile of the final cluster to which it has contributed (see cluster 2 on Figure 5). The rainfalls of the 08th of January 2018 were locally considerable (more than 60 mm, see Figure 7-a) whereas the process classifies this day in a finally removed cluster due to too few rainfalls on average.

We now present whether the final result of the process is valid in correctly classifying these days and decide which weather trends affect them. We determine that by comparing the profiles of these days (the one visible on Figure 7-b, obtained after applying equations 1 and 2 on their raw data) with those of

the final clusters (FC), which are typical of each detected weather trend. We calculate the euclidean distances between all profiles. The day to analyze is then associated with the closest final cluster (i.e., to the most similar fictive day given by the multilevel clustering process) and corresponds to the respective trend.

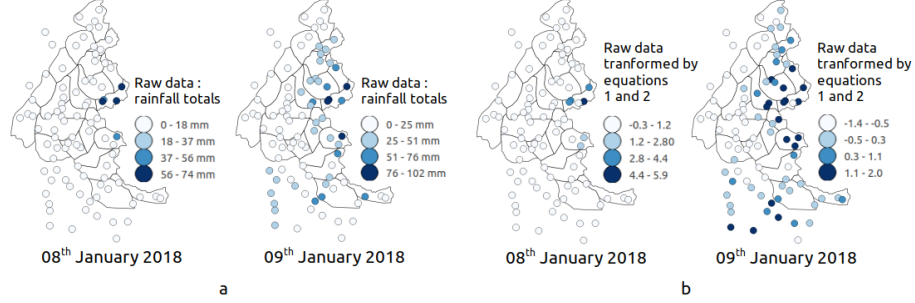


Fig. 7: Rainfall totals and profiles of the 08th and 09th January 2018

	FC 1	FC 2	FC 3	FC 4
<b>08th Jan. 2018</b>	13.71	13.94	9.97	15.07
<b>09th Jan. 2018</b>	12.82	13.94	8.30	15.50

Fig. 8: Distances between days to analyze and final clusters

According to the result (Figure 8), both days are affected to the resulting cluster number 3, i.e., the one corresponding to a flux coming back from the east. This conclusion agrees with our first assumption and the expert knowledge we have for these days. Thus, the result of the process is valid as it permits correctly classifying the days. It corrects their cluster, even if they contribute wrongly, during the process. That shows the weight of the global and the validity of the final result given by the multilevel clustering process.

The four detected trends output by the multilevel clustering process have been validated based on expert knowledge. Indeed, we cannot establish a correct classification rate since the expected cluster for each day is not known in advance (no complete ground truth on the dataset). In addition, the cluster given to each day individually does not matter so much as the clustering process is useful to finally output and detect the main trends affecting the studied area. Thus, the multilevel clustering process highlights the different areas impacted by precipitation according to different trends. This result is particularly useful in avalanche danger management as these areas fit the ones impacted by spontaneous avalanches that occurred on the corresponding days. Finally, we will massively test the framework and do some further validation on future data (recorded for the following winters).

## 4 Discussion

First, the multilevel clustering process can be adapted and used with other clustering methods than those used here. Depending on the data, some clustering

methods are more suitable than others. Depending on the expected shape of the clusters, it will be, for example, more appropriate to use K-means (circular) or GMM (ellipsoidal) [20], or for any other reason, another clustering method.

In our work, we tried with AP or GMM instead of K-means in the first level clustering. But, these two methods are not only a consolidation of HC results, as was the case for K-means. Then, the sequence of methods was less fluent with AP: the parameter of preference has to be adapted to obtain the number of clusters given by HC, which needs several tries. Nevertheless, the result obtained at the end was very close. With GMM instead of K-means, results were only slightly different. But as we did not dispose of adequate ground truth (only a few days among 1455 available), we did not have any possibility to be sure that one result was better than another. In addition, GMM requires adapting some parameters to obtain the best possible result (most relevant with the data) and, without any ground truth, it was impossible to estimate what changes were needed. We can not derive conclusions on any improvement or degradation of the result. With the constraints to work without any ground truth and any expert intervention, we prefer keep using the K-means method which was a consolidation of HC and does not need any other parameter to be set and refined.

A possible follow up to this work is the refinement of the profiles of the four obtained fictitious days, corresponding to typical representants of the four detected weather trends. The validation step showed that days, which the process misclassifies, were finally well estimated based on their distances to the four final fictitious days (see Section 3.2). This way, by recalculating to which final cluster each day should belong to, each cluster's profile can be updated. Thus, the updated result will be based on well-classified individuals, which should improve the correctness and precision. Thus it will be a better base to estimate, for each weather trend, which area encounters similar rainfall conditions.

Then, what should follow this work is to study, for each weather trend, which stations have the same features. That means studying which stations receive similar amounts of new rainfalls according to different weather trends. Even if we already get a rough idea of these zoning through the map representing the profiles (see Figure 5), some additional work is required to precise which areas are similar in terms of rainfall conditions, for each trend. It will take place in the general objective firstly mentioned in the introduction and the specific context of our work, to finally link areas impacted by rainfalls and avalanche observations for given days.

## 5 Conclusion

Our multilevel clustering process is based on two clustering levels that use classical methods in sequence: 1) PCA  $\rightarrow$  HC  $\rightarrow$  K-means, and 2) AP. Clustering in two levels masks the misclassification of some individuals by giving a correct global result. At the end, the process gives four typical profiles, which are equivalent to fictitious days presenting the features of the detected weather trends.

The first level clustering does not always detect all the trends present in the data subset. That is why, to increase the capacity to detect all the possible main

trends at least once, we increase the number of studied days in the input of the process. We divide these days into several data subsets. Otherwise, the obtained result does not correspond to what we are looking for: it classifies days according to intensity instead of the location of the rainfalls.

On the other hand, the first level clustering sometimes gives several clusters that correspond to only one trend (i.e., profiles with the same shape but differences in values), whereas we want to detect a single final cluster for each weather trend. We do not stop only after this first level, to aggregate the similar clusters. At the end of the second level clustering, we obtain only one final cluster for each detected weather trend. The role of this second level is also to merge the results obtained due to the  $n$  iterations of the first level clustering.

Moreover, if some days could be misclassified, the global result is in adequacy with what is intended. It is possible to extract the associated area affected by precipitations. That is why the principle of our process is not to check precisely in which cluster each day is affected but to consider only the global result obtained thanks to means and two-level clustering.

The strength of our process is that despite some misclassifications on initial individuals (i.e., days), the weight of the collective gives a satisfying final result. Indeed, when we use the result to estimate which weather trend a day should be associated with, the estimation is correct. By associating a day to the closest fictitious day among those given in result by our multilevel clustering process, we can affect the right trend to the analyzed day. That means the output of our process is good enough to correctly decide: which weather trend affects a day among the main detected trends.

In the future, the developed process will be integrated into the decision support method CRISTAL [2]. As the process automatically assesses the extent of the area which makes sense for a given avalanche observation, it will fill one of the six parameters on which CRISTAL relies. Indeed, as the multilevel clustering result gives a better understanding of the areas affected by precipitation under different meteorological trends, it makes it possible to specify a criterion that is vague when it is taken into account in the existing decision-making frameworks (as mentioned, for example, in the Obvious Clues Method). Thanks to this result, we can better formalize the boundaries of the area subject to similar danger as the one that led to avalanching.

**Acknowledgements.** The CIME project is supported by the European cross-border cooperation program Interreg France-Switzerland 2014-2020 and has been awarded a European grant (European Regional Development Fund) covering 60% of the total French cost.

We thank the snow expert Alain Duclos, who provided us the information on some typical days representative of the trends, especially concerning the third detected one.

## References

1. M.J. Casado, M.A. Pastor, F.J. Doblas-Reyes: Links between circulation types and precipitation over Spain, *Physics and Chemistry of the Earth* 35, pp. 437-447, 2010.
2. A. Duclos: *Nivologie pratique : les 4 modes de vigilance encadrée. Neige et avalanches* 160, pp. 7-9, 2018.

3. B. J. Frey, D. Dueck: Clustering by Passing Messages Between Data Points, *Science*, vol. 315, no. 5814, pp. 972-976, Feb. 2007.
4. M. Huan, R. Lin, S. Uang, T. Xing: A novel approach for precipitation forecast via improved K-nearest neighbor algorithm, *Advanced Engineering Informatics* 33, pp. 89-95, 2017.
5. M. Irannezhad, A.K. Ronkanen, S. Kiani, D. Chen, B. Klove: Long-term variability and trends in annual snowfall/total precipitation ratio in Finland and the role of atmospheric circulation patterns, *Cold Regions Science and Technology* 143, pp. 23-31, 2017.
6. M. Landrø, A. Hetland, R. Engeset, G. Pfuhl: Avalanche decision-making frameworks: Factors and methods used by experts. *Cold Regions Science and Technology* 170, 2020.
7. M. Landrø, G. Pfuhl, R. Engeset, M. Jackson, A. Hetland: Avalanche decision-making frameworks: Classification and description of underlying factors. *Cold Regions Science and Technology* 169, 2020.
8. M. Lemus-Canovas, J. A. Lopez-Bustins, L. Trapero, J. Martin-Vide: Combining circulation weather types and daily precipitation modelling to derive climatic precipitation regions in the Pyrenees, *Atmospheric Research* 220, pp. 181-193, 2019.
9. F. Louchet: *Snow Avalanches: Beliefs, Facts, and Science*. 112 p., 2021. ISBN: 9780198866930
10. I. McCammon, 2006. Obvious Clues Method: A Users Guide. *The Avalanche Review* 25 (2), 2006.
11. J. Meng, H. Hao, Y. Luan: Classifier ensemble selection based on affinity propagation clustering, *Journal of Biomedical Informatics*, 60, pp.234-242, 2016.
12. F. Murtagh, P. Legendre: Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion?, *Journal of Classification*, 31, p.274-295, 2014. DOI: 10.1007/s00357-014-9161-z
13. F. Pagnier, D. Coquin, F. Pourraz, H. Verjus, G. Mauris: Classification des précipitations sur les massifs alpins français. ORASIS 2021, Centre National de la Recherche Scientifique [CNRS], Saint Ferréol, France, Sep 2021. hal-03339626
14. J.P. Praene, B. Malet-Damour, M.H. Radanielina, L. Fontaine, G. Rivière: GIS-based approach to identify climatic zoning: A hierarchical clustering on principal component analysis, *Building and Environment* 164, 2019.
15. M. B. Richman, I. Adrianto: Classification and regionalization through kernel principal component analysis, *Physics and Chemistry of the Earth* 35, pp. 316-328, 2010.
16. J. Schweizer, K. Kronholm, J.B. Jamieson, K.W. Birkeland: Review of spatial variability of snowpack properties and its importance for avalanche formation. *Cold Regions Science and Technology* 51, pp. 253-272, 2008.
17. S. Tufféry: *Data mining et statistique décisionnelle : L'intelligence des données*, 4ème édition, Editions Technip, Paris, 2012.
18. J.H. Ward: Hierarchical Grouping to Optimize an Objective Function. *Journal of the American Statistical Association*, 1963. <https://doi.org/10.1080/01621459.1963.10500845>
19. B. Zahraie, A. Rooszbahani: SST clustering for winter precipitation prediction in southeast of Iran: comparison between modified K-means and genetic algorithm-based clustering methods. *Expert Systems with Application* 38, pp. 5919-5929, 2011.
20. L. Zhao, Z. Shang, J. Tan, X. Luo, T. Zhang, Y. Wei, Y. Yan Tang: Adaptive parameter estimation of GMM and its application in clustering, *Future Generation Computer Systems*, Volume 106, pp. 250-259, 2020.