

## Automated detection of case clusters of waterborne acute gastroenteritis from health insurance data – pilot study in three French districts

Loïc Rambaud, Catherine Galey and Pascal Beaudeau

### ABSTRACT

This pilot study was conducted to assess the utility of using a health insurance database for the automated detection of waterborne outbreaks of acute gastroenteritis (AGE). The weekly number of AGE cases for which the patient consulted a doctor (cAGE) was derived from this database for 1,543 towns in three French districts during the 2009–2012 period. The method we used is based on a spatial comparison of incidence rates and of their time trends between the target town and the district. Each municipality was tested, week by week, for the entire study period. Overall, 193 clusters were identified, 10% of the municipalities were involved in at least one cluster and less than 2% in several. We can infer that nationwide more than 1,000 clusters involving 30,000 cases of cAGE each year may be linked to tap water. The clusters discovered with this automated detection system will be reported to local operators for investigation of the situations at highest risk. This method will be compared with others before automated detection is implemented on a national level.

**Key words** | automated detection, drinking water, France, gastroenteritis, health insurance data, outbreak

Loïc Rambaud (corresponding author)

Catherine Galey

Pascal Beaudeau

French Institute for Public Health Surveillance,

12 rue du Val d'Osne,

94415 Saint Maurice,

France

E-mail: l.rambaud@invs.sante.fr

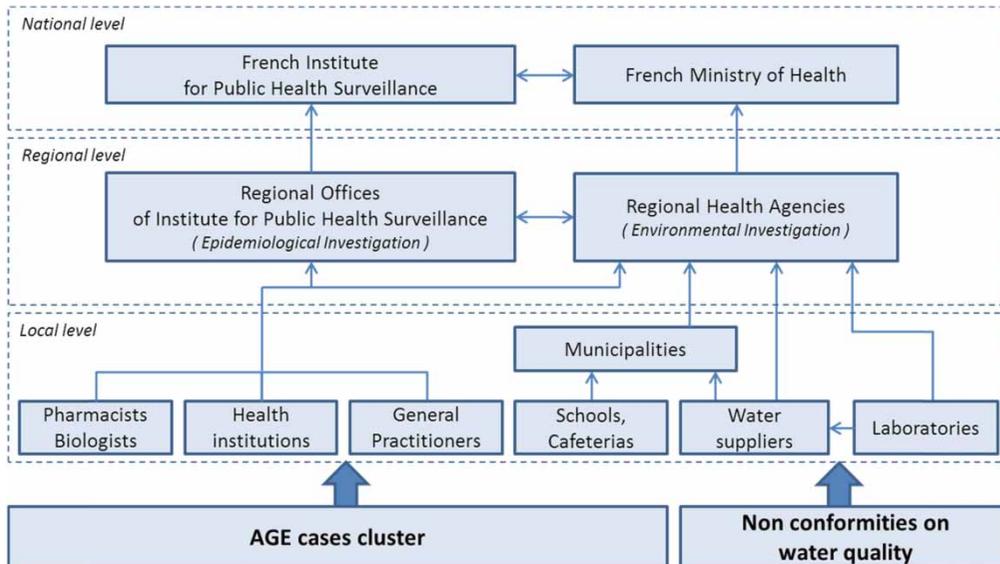
### INTRODUCTION

Surveillance of risks of fecal origin linked to the drinking water supply traditionally focuses on outbreaks (Brunkard *et al.* 2011). Since the early 2000s, the French Institute for Public Health Surveillance (InVS) has identified one to two waterborne outbreaks of acute gastroenteritis (AGE) (Beau-deau *et al.* 2008). Their detection is based on reporting by local officials, operators or physicians of a currently happening cluster of AGE cases, or from laboratory results of water analysis exceeding regulatory levels (Figure 1). These situations are supported by local health authorities (Regional Health Agencies and regional offices of the French Institute for Public Health Surveillance), and usually need a laboratory follow up to confirm the waterborne origin of the outbreak and microbiological infection of the cases. Small-scale clusters, difficult to identify in the field, probably go unnoticed.

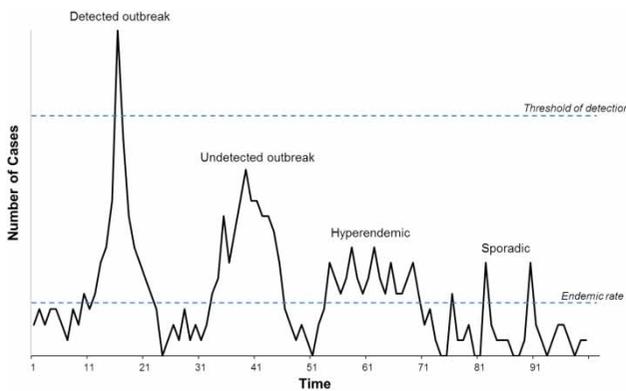
Improvement in the surveillance of the risk of waterborne diseases thus requires an increased capacity to detect AGE

case clusters (Figure 2). The development of medical-administrative databases provides an opportunity to do so. Since 2009, Sniiram (a national medical information database) (Tuppin *et al.* 2010) has been used to produce a daily estimate of the number of AGE cases seen by doctors (cAGE) (Bounoure *et al.* 2011) throughout France, at the scale of the municipality. InVS already uses this indicator to validate and retrospectively characterize outbreaks reported or suspected by local health authorities (Rambaud *et al.* 2011).

This new approach to the surveillance of waterborne disease risks in France uses Sniiram to search routinely for cAGE case clusters. Repeated at regular intervals (for example, every 2 months), it enables the retrospective identification of the most vulnerable water supply systems. The detection of case clusters is based on the demonstration of an abnormally high number of observed cAGE cases (NOC) in the spatiotemporal unit considered (e.g.,



**Figure 1** | French system for surveillance of waterborne outbreaks.



**Figure 2** | Limits to outbreak detection (adapted from Frost *et al.* 1996). Reprinted from *Journal AWWA* 88 (9), by permission. Copyright 1996 American Water Works Association.

municipality per week) compared with the number of expected cases (NEC). Several methods exist for this purpose. Although the most common methods (e.g., control charts, historical averages) do not appear appropriate to the spatialized nature of the data, Kulldorff's space-time scan (Kulldorff *et al.* 2005) makes it possible to take the spatial structure of water supply networks into account. Nonetheless, its implementation as part of a waterborne disease cluster detection system would require a thoroughly accurate mapping of these networks, which is currently unavailable at the national level. A simple, pragmatic method

inspired by practice in the field has been developed and was assessed during the pilot study that is described here. The prospects offered by the detection of cAGE clusters for the ultimate goal of preventing waterborne disease outbreaks will then be discussed.

## MATERIAL AND METHODS

### Health data

The method of identifying cAGE cases from Sniiram data was described earlier (Bounoure *et al.* 2011). The Sniiram database covers the French population almost exhaustively and makes it possible to obtain individual medical consultation data on a same-day basis, localized to the patient's town of residence. Cases with consultations more than 50 km from their home are excluded. A survey of the customers of several pharmacies (Bounoure *et al.* 2011) estimated this method's sensitivity and specificity at 0.89, compared with cAGE cases seen by a doctor certifying the diagnosis according to the consensus syndrome definition (Majowicz *et al.* 2008).

Cumulative cases per town are tallied weekly (Monday to Sunday), before the search for clusters. A weekly observation period fits the duration of waterborne disease outbreaks (1–3 weeks) better and avoid the effects

associated with daily variations due to fluctuations in work, such as Sunday closings.

### Study period and area

The study period ran from Monday January 5, 2009, to Sunday December 30, 2012. Data are missing for the period from March 21 to May 1, 2011, inclusive.

The study sector includes three French districts, Puy-de-Dôme, Isère, and Gironde (Figure 3), representative of the country's diversity (city vs. countryside, plains vs. mountains). According to the National Institute of Statistics and Economic Studies (INSEE), the total population of the study area in 2009 was 3.215 million inhabitants. The distribution of towns by population size is similar in Isère and Gironde. The Puy-de-Dôme has fewer towns, and a higher percentage of them have small populations (Table 1). The

number of events involving accidental microbiological pollution of the water supply in these sectors is close to the national mean, and did not differ significantly between these districts (Beaudeau 2010). The distributed water is mainly of underground origin in the three districts, as 99.1% of the water supply network use this type of resource in Isère, compared with 99.2% in Gironde and 98.1% in Puy-de-Dôme. Distributed water meets the microbiological standards for more than 95% of the supplied population in Isère and Gironde, and for more than 90% in Puy-de-Dôme.

### Detection of cAGE case clusters

We developed two methods, named A and B, both based on the comparison of the town tested and a reference area. For each week and each town, weekly cluster detection was performed in two stages: (1) determination of the NEC, by steps



Source :  
©IGN-GéoFLA®, 2011 ;  
IGN-Géoportail, Elevation slope  
InVS, juillet 2015

Figure 3 | French districts included in the study area.

**Table 1** | Distribution of towns in the study area by population (INSEE data, 2009)

| Town size (inhabitants) | Gironde |      | Isère |      | Puy-de-Dôme |      | Total |      |
|-------------------------|---------|------|-------|------|-------------|------|-------|------|
| (0, 100)                | 14      | 3%   | 23    | 4%   | 29          | 6%   | 66    | 4%   |
| (100, 500)              | 210     | 39%  | 142   | 27%  | 238         | 51%  | 590   | 38%  |
| (500, 2,000)            | 202     | 37%  | 243   | 46%  | 154         | 33%  | 599   | 39%  |
| (2,000, 10,000)         | 94      | 17%  | 109   | 21%  | 39          | 8%   | 242   | 16%  |
| (10,000, 50,000)        | 19      | 4%   | 13    | 2%   | 9           | 2%   | 41    | 3%   |
| (50,000 et +)           | 3       | 1%   | 1     | 0%   | 1           | 0%   | 5     | 0%   |
| Total                   | 542     | 100% | 531   | 100% | 470         | 100% | 1,543 | 100% |

appropriate to each method (Figure 4) and (2) a common test comparing the NOC and the NEC.

### Determination of the NEC

Regardless of the method, A or B, the entire district to which the town tested belongs is used as the reference area.

Method A (Figure 4) is based on a spatial comparison of the incidence of cAGE cases between the town tested and the reference area. The NEC is the product of the median incidence rate in the reference area multiplied by the size of the population in the municipality tested. It is preferable to use the median because it, unlike the mean, is less sensitive to extreme (e.g., associated with outbursts) or outlying values. This reference median rate can nonetheless be zero, especially in the summer in rural districts, where most towns are small; to eliminate this disadvantage, towns with fewer than 500 inhabitants are excluded from its calculation.

Method B (Figure 4) is based on a temporal change comparison of the incidence of cAGE cases between the town tested and the reference area. The relative variation in incidence rates is calculated for each town in the reference area as the ratio between the cAGE incidence rate during the target period and the cAGE incidence rate for a control period. The relative variation of the reference area corresponds to the median of this set of values (again excluding towns with fewer than 500 inhabitants). Then, to produce the NEC, this relative variation is multiplied by the mean weekly incidence rate observed during the control period for the town tested. The control period chosen covers the 4 weeks between the fifth ( $W_{-5}$ ) and second weeks ( $W_{-2}$ ) before the target week ( $W_0$ ). A one-week buffer is thus placed between the target week and the control period, so that the latter does not include the beginning of a potential outbreak.

In low incidence periods (spring, summer), it is possible to observe no cases in the smallest towns during the control period. In these situations, we imputed the value 1 to the number of cases observed in the town during the control period to allow the calculation of variations in the town incidence rates. This device tends to reduce the number of clusters detected by method B.

### Test of NOC versus NEC

The detection of suspected weekly clusters is based on the combined use of three criteria:

1. **A relative risk (RR) of cAGE  $\geq 2$ .** The RR is estimated by the ratio of the NOC to the NEC. This choice is based on the observation that a RR in the order of 2 can be observed without the presence of any environmental risk factor between populations that differ solely by their age distribution (Tam et al. 2003).
2. **A health impact  $> 5$  cases.** The health impact is defined as the excess number of observed compared to expected cases (NOC-NEC). This choice is intended to exclude foodborne clusters or those due to personal contact, both of which are most frequent among small clusters (Delmas et al. 2006).
3. **A probability  $p < 1.10^{-5}$ .**  $P$  corresponds to the significance of the statistical test performed. The null hypothesis is that the NOC of the town tested during the target week is equal to the NEC. An exact test is used, and we assume that the NOC has a Poisson distribution with parameter NEC. In view of the repetition and non-independence of the tests, the p-value estimate is biased by default but preserves the relative order of significance.

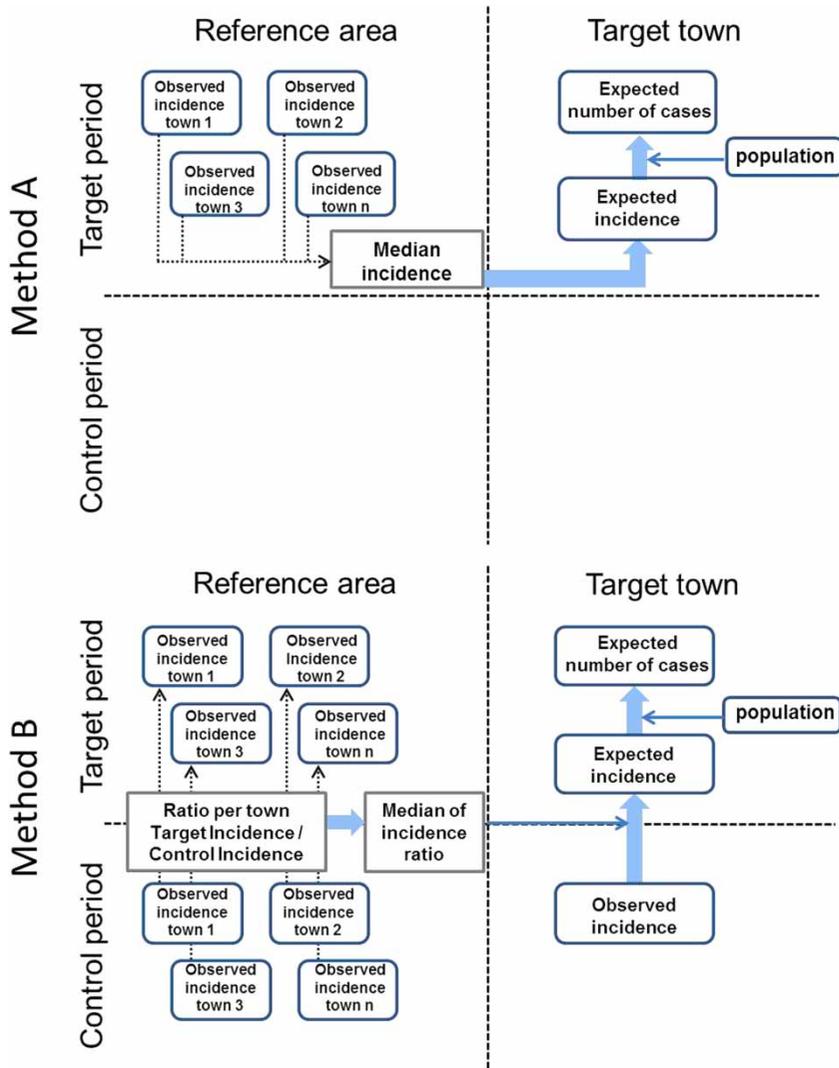


Figure 4 | Principle for calculating the NEC.

### Cluster consolidation

The weekly clusters detected simultaneously by method A (list A) and method B (list B) are selected. They are considered to have been detected by methods  $A \cap B$  (list  $A \cap B$ ). In each of the three lists, weekly clusters of the same towns that are consecutive in time are combined into a single cluster, called a consolidated cluster, with a duration of a week or more. The impact and the RR of a consolidated cluster are recalculated from the cumulative NOCs and NECs of the weekly clusters included in the consolidated cluster.

## RESULTS

### Frequency of clusters

Our systematic exploration detected 826 weekly clusters by method A and 543 by method B; 210 were identified by both methods (Figure 5). The detection rate was 6.5/10,000 municipalities  $\times$  weeks ( $M \times W$ ) and varied between 5/10,000  $M \times W$  in Isère and 8/10,000  $M \times W$  in Puy-de-Dôme. Over the entire study area, method A detected 1.5 times more weekly clusters than method B. This rate varied between 0.8 and 2 according to district. The rule requiring

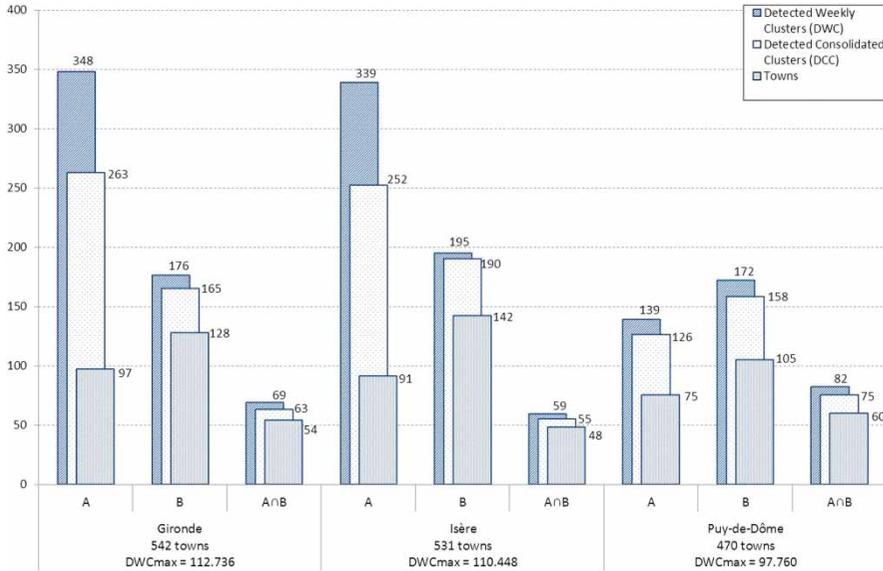


Figure 5 | Detection of cAGE case clusters and towns concerned according to the method used (Puy-de-Dôme, Isère and Gironde; 2009–2012).

intersection ( $A \cap B$ ) selects fairly drastically, as only 18% of the weekly clusters detected by either method ( $A \cup B$ ) are found in the intersection list  $A \cap B$ . This selection effect varies between 12 and 36%, according to district.

The 210 weekly clusters detected in the complete study area by methods  $A \cap B$  were consolidated into 193 clusters.

### Distribution of the clusters detected by size of municipality

More than two thirds of the consolidated clusters detected by methods  $A \cap B$  were concentrated in municipalities of

500 to 10,000 inhabitants (55% of the towns) and barely 2% involve towns larger than 50,000 inhabitants (0.3% of all towns) (Table 2). The frequency of detection of consolidated clusters by methods  $A \cap B$  is zero for towns with fewer than 100 inhabitants, and increases with town population up to 50,000 inhabitants. For towns with fewer than 100 inhabitants, only method A detected any clusters. For each method, the distribution of the number of consolidated clusters detected differed by town population-size groups; this difference was substantial for the group of towns with 10,000–50,000 inhabitants, for which method A detected more clusters.

Table 2 | Distribution of consolidated clusters detected, according to the population of the affected town (Puy-de-Dôme, Isère and Gironde; 2009–2012)

| Town size (inhabitants) | List A        |       | List B        |       | List $A \cap B$ |       | Number of towns |                           |        |
|-------------------------|---------------|-------|---------------|-------|-----------------|-------|-----------------|---------------------------|--------|
|                         | Number of DCC | (%)   | Number of DCC | (%)   | Number of DCC   | (%)   | Total           | Concerned by $\geq 1$ DCC | (%)    |
| (0, 100)                | 1             | (0)   | 0             | (0)   | 0               | (0)   | 66              | 0                         | (0.0)  |
| (100, 500)              | 81            | (12)  | 47            | (9)   | 30              | (15)  | 590             | 27                        | (4.6)  |
| (500, 2,000)            | 192           | (30)  | 238           | (46)  | 70              | (36)  | 599             | 58                        | (9.7)  |
| (2,000, 10,000)         | 171           | (27)  | 182           | (36)  | 67              | (35)  | 242             | 60                        | (24.8) |
| (10,000, 50,000)        | 192           | (30)  | 41            | (8)   | 23              | (12)  | 41              | 16                        | (39.0) |
| (50,000 et +)           | 4             | (1)   | 5             | (1)   | 3               | (2)   | 5               | 1                         | (20.0) |
| Total                   | 641           | (100) | 513           | (100) | 193             | (100) | 1,543           | 162                       | (10.5) |

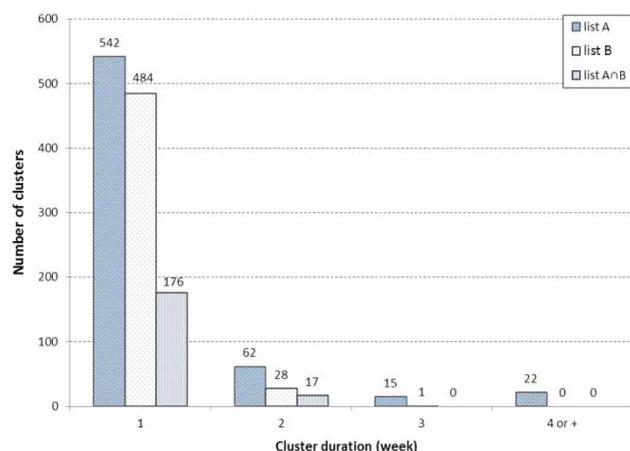
DCC = detected consolidated cluster.

## Duration of clusters

More than 90% of the clusters on list  $A \cap B$  (176) lasted only 1 week (Figure 6), and none more than 2 weeks. Method B essentially detected clusters of 1 or 2 weeks (only 1 consolidated cluster of 3 weeks or longer), while only method A detected clusters lasting longer than 3 weeks (3.4%), including one that lasted 11 weeks.

## Municipalities at chronic risk

Between 2009 and 2012, slightly more than 10% ( $n = 193$ ) of towns in the study area were involved in at least one consolidated cluster (Table 3). Repeated clusters occurred in only 14% of the positive towns (with at least one cluster during the 4-year study period), and in only 1.5% of all towns.



**Figure 6** | Distribution of consolidated clusters of cAGE cases according to their duration (Puy-de-Dôme, Isère and Gironde; 2009–2012).

**Table 3** | Distribution of cluster-positive towns according to the number of consolidated clusters detected (Puy-de-Dôme, Isère and Gironde; 2009–2012, N towns = 1,543)

| Number of DCC by town | List A          |       | List B          |       | List A ∩ B      |       |
|-----------------------|-----------------|-------|-----------------|-------|-----------------|-------|
|                       | Number of towns | (%)   | Number of towns | (%)   | Number of towns | (%)   |
| 1                     | 164             | (62)  | 270             | (72)  | 139             | (86)  |
| 2                     | 47              | (18)  | 79              | (21)  | 16              | (10)  |
| 3–5                   | 30              | (11)  | 26              | (7)   | 7               | (4)   |
| 6 or +                | 22              | (8)   | 0               | (0)   | 0               | (0)   |
| Total                 | 263             | (100) | 375             | (100) | 162             | (100) |

DCC = detected consolidated cluster.

The maximum number of repetitions was four consolidated clusters for one town.

More consolidated clusters were detected by method A but they concerned fewer municipalities than those detected by method B (263 municipalities for A *vs.* 375 for B). Thus, repetitions of consolidated clusters in the same town were more easily detected by Method A. It was the only method to show repetitions of six consolidated clusters or more by town; the maximum was 32, compared with only five per municipality for method B.

## Number of cases involved

Of the consolidated clusters detected on list  $A \cap B$ , 40% involved fewer than 10 cases of cAGE each (Table 4) and accounted for 16% of all cases of detected cAGE. On the other hand, 6% of the consolidated clusters involved 50 cases or more and accounted for 28% of all the cAGE cases detected. The largest cluster detected involved 199 cases over a 2-week period.

Method B tended to detect smaller consolidated clusters more easily. Accordingly, 58% of the clusters detected by method B included fewer than 10 cases compared with only 29% for method A, accounting for 31% and 8%, respectively, of all the clustered cases. More than two thirds of the cases detected by method B came from clusters of fewer than 20 cases compared with 25% for method A (Table 4).

## Seasonal variations in detection

The mean number of monthly consolidated clusters detected by method  $A \cap B$  was 16. It ranged from 6 in May to 35 in December (Figure 7). Regardless of the method used, twice as many clusters were detected in December. Seasonality did not seem to differ between any of the three methods A, B, and  $A \cap B$ .

## DISCUSSION

### In search of specificity for waterborne cAGE cases clusters: a public health objective

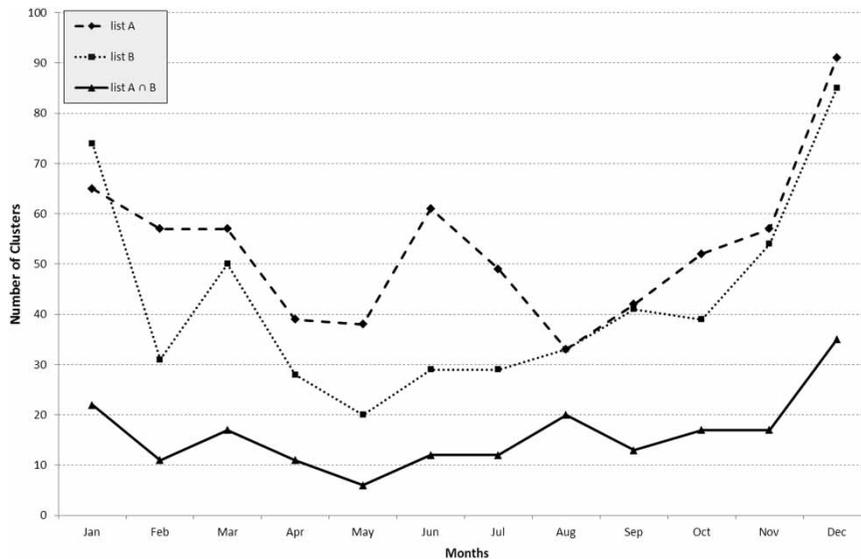
Supposing that the frequency of cluster detection in the three pilot districts is representative of the national situation, we

**Table 4** | Distribution of consolidated clusters detected according to method A, B or A ∩ B and the number of cAGE cases involved (Puy-de-Dôme, Isère and Gironde; 2009–2012)

| Health impact (cases by cluster) | List A |       |        |       | List B |       |       |       | List A ∩ B |       |       |       |       |       |
|----------------------------------|--------|-------|--------|-------|--------|-------|-------|-------|------------|-------|-------|-------|-------|-------|
|                                  | NClu   | (%)   | Nca    | (%)   | NClu   | (%)   | Nca   | (%)   | NClu       | (%)   | Nca   | (%)   | Nca*  | (%)   |
| (6, 10)                          | 186    | (29)  | 1,342  | (8)   | 299    | (58)  | 2,084 | (31)  | 79         | (41)  | 533   | (16)  | 545   | (13)  |
| (10, 20)                         | 212    | (33)  | 3,004  | (18)  | 146    | (28)  | 2,008 | (30)  | 68         | (35)  | 932   | (28)  | 872   | (20)  |
| (20, 50)                         | 179    | (28)  | 5,133  | (31)  | 54     | (11)  | 1,436 | (21)  | 35         | (18)  | 927   | (28)  | 1,165 | (27)  |
| (50, 100)                        | 44     | (7)   | 2,918  | (18)  | 10     | (2)   | 676   | (10)  | 9          | (5)   | 584   | (18)  | 495   | (12)  |
| [100 or +                        | 20     | (3)   | 4,253  | (26)  | 4      | (1)   | 598   | (9)   | 2          | (1)   | 356   | (11)  | 1,194 | (28)  |
| Total                            | 641    | (100) | 16,650 | (100) | 513    | (100) | 6,802 | (100) | 193        | (100) | 3,332 | (100) | 4,271 | (100) |

NClu = number of clusters; Nca = number of cases of cAGE.

\*Calculated from health impacts obtained by method A.

**Figure 7** | Monthly distribution of the total number of clusters of AGE cases detected (Puy-de-Dôme, Isère, and Gironde; 2009–2012).

have estimated that between 1,000 and 2,000 clusters could be detected annually in France. These estimates are orders of magnitude to be validated and specified in the future. The automated detection of case clusters nonetheless represents crucial progress in terms of sensitivity compared with the traditional method based on reporting by local operators.

These estimates exceed the investigative capacity of the local authorities responsible for controlling the quality of the water supply. Implementation of the method in the field thus depends on improving its specificity, that is, of separating the waterborne clusters from those which we do not wish to detect, that is: clusters due to chance, with no identifiable cause; false clusters associated with methodological biases;

clusters propagated by human-to-human transmission and part of the winter virus epidemic; foodborne outbreaks; and clusters associated with swimming in pools or in recreational water.

It can be difficult to distinguish a waterborne cAGE outbreak from another cAGE outbreak by its epidemic curve alone. Additional information is thus necessary to support its waterborne origin.

One useful type of information for this purpose is knowledge of the geographic contours of water supply networks, which do not necessarily correspond to those of municipalities. Specifically, the sharing of a network by several municipalities with outbreaks suggests that it plays a role

in the outbreak onset. Nonetheless, only 20% of French towns share a water supply network with another town. Moreover, the national map of water-supply networks is not yet exhaustive and does not list temporary connections, which are frequent in tourist sectors. Associating a drinking water distribution area and an outbreak on one or several municipalities therefore requires systematic validation in the field.

The list of detected cAGE consolidated clusters will be transmitted to local health authorities so that they can conduct environmental investigations to determine if clusters have waterborne origins, identify their causes, and take appropriate prevention measures. The transmission of the information collected in these field surveys to a national level and their analysis should also be promoted. These data will make up a national surveillance system for waterborne risk factors and will thus help to direct prevention of such outbreaks through the development of relevant regulations and the promotion of good practices.

The method of cluster detection must privilege specificity over sensitivity to avoid the dissipation of limited field resources in the investigation of non-waterborne clusters. Accordingly, the list transmitted to local health authorities will be sorted by the number of cAGE cases involved in the clusters detected. For example, clusters of more than 20 cAGE cases account for only 24% of the clusters found by method  $A \cap B$  (but 67% of those counted with method A, see below).

Field studies could be coordinated with the implementation of the water safety plans recommended by WHO (WHO 2005), especially in cases of substantial or repeated outbreaks in the same town: a second list by municipality could be created and sorted by the number of cumulative clustered cAGE cases in recent years (e.g., 2009), to guide the deployment of water safety plans toward the towns most frequently affected.

Non-waterborne clusters can present particularities that make it possible to identify them without a field investigation. Family foodborne outbreaks involve a small number of cases (Delmas *et al.* 2006) and will not appear among the clusters detected, and still less among the clusters investigated in the field. Some of the foodborne disease outbreaks occurring in institutional food services will have already been reported and investigated and will therefore

not be candidates for a new investigation. Those not previously identified may subsequently be found by their age structure (for example, ratio of the numbers of children and adults), which may be imbalanced compared to that of waterborne clusters, which affect all age groups. Foodborne disease outbreaks occurring in the cafeterias of schools, companies, or retirement homes can thus be excluded *a priori* from the investigations. Foodborne disease outbreaks reported and associated with festive or community meals can lead to confusion and unnecessary field surveys.

Among the cAGE cases clusters associated with swimming, those concerning concentrated tourist populations (campgrounds, resorts) are not included in the surveillance. Confusion with an outbreak linked to the contamination of a town pool or other recreational swimming site cannot be ruled out (Smith *et al.* 2006).

Information from the field remains predominant in assessing the likelihood of the waterborne origin of a cAGE cluster detected. The attribution of a waterborne origin to an outbreak traditionally requires evidence – either microbiological (pathogen strain identified simultaneously in water and in patients' stools), or epidemiological (analytic studies showing a significant association between the quantity of water consumed and the risk of AGE and allowing no alternative explanations) (Tillet *et al.* 1998). For the surveillance system proposed, collecting adequate water and stool samples will be impossible because of the delay between the occurrence and the detection of the cluster. In this case, epidemiologic evidence of waterborne origin will not be furnished by a descriptive epidemiologic analysis alone. In a first step, evidence to attribute the waterborne origin will come from further analysis of detected clusters, for example, the size of the cluster, its age structure and its spatial extent (evidence limited to cases of networks serving several towns). In a second step, an environmental survey conducted in partnership with drinking water operators can lead to the identification of unfavorable conditions for safe water production.

The environmental survey thus aims essentially to show events during the pre-outbreak period likely to indicate a waterborne source: unfavorable weather or water conditions (intense rain, rising water levels, floods, pollution, overflows from an upstream sewage system, etc.); dysfunction in the

water production system (a breakdown in disinfection, defective clarification, work on the network, etc.); an increase in user complaints (odor, taste, color suggesting organic pollution); and microbiological analyses of the water showing that it fails to comply with regulatory standards.

The concomitance of one of these events and of a cluster of cAGE cases supports the hypothesis of a waterborne origin of the cluster.

### In search of specificity: methodological aspects

The criteria used to detect clusters ( $RR \geq 2$ , impact  $>5$  and  $p < 1.10^{-5}$ ) by method A, on the one hand, and by method B on the other, are the first barrier against false positives. But it is the selection of clusters found simultaneously by A and B (methods  $A \cap B$ ) that is the key factor for specificity: only 18% of the clusters found by A or B are common to both. That is, this selection makes it possible to eliminate the false positives specific to either method A or method B.

If we assume a Poisson distribution of the weekly counts of cases common to both methods, the multiplication of statistical tests performed engenders for each district in the study area one false positive due to chance every 4 years, with method A and the same number with method B. Although this probability is underestimated, it is improbable that the intersection process would finally retain one of these false positives.

The selection condition that  $RR \geq 2$  is intended to limit the false positives inherent in the use of method A and attributable to the AGE incidence rate variations between towns, induced by demographic or socio-economic factors.

The demographic bias probably induces the greatest bias in the comparison of municipal incidence rates. Two phenomena combine to generate this bias: (1) the variability in the age structure of the populations in different towns; in towns with more than 500 inhabitants, the ratio of the number of adults to the number of children thus varies from 2 to 10; and (2) the cAGE incidence rate is four times higher in children than in adults.

The demographic bias is clearly expressed in municipalities with a population between 10,000 and 50,000 inhabitants (Table 2). In this stratum, the ratio of the numbers of children to the number of adults is 0.30 in the

municipalities for which a cluster is detected by method A compared with 0.21 for the others; this difference suggests overdetection by method A. Age standardization for the number of cAGE cases should be considered if method A is to be used alone; this nonetheless raises a size or power problem (smaller weekly numbers of cases in smaller towns). This flaw is corrected here by method B, which is insensitive to this bias because it is based on the comparison of relative variations in incidence between the target town and the district. The same is true for socioeconomic biases, such as access to physicians or educational level, which may condition these consultations (Tam *et al.* 2003).

The intersection of methods A and B can also produce a loss of sensitivity; nonetheless, this is not expressed by false negatives, but by truncating the clusters. Beyond the second week of an outbreak, the control period used by method B mechanically includes the beginning of the outbreak, which thus reduces the probability of detection of a weekly cluster during the third week or later, that is, artificially reduces the duration and impact of the consolidated clusters detected. To compensate for this defect, we propose to attribute the impact assessed by method A to the clusters detected by methods  $A \cap B$ . This triples the estimate of the cumulative number of cases included in the clusters of more than 100 cases (Table 4). The clusters of a duration of three weeks or more thus rise from 0 to 2.5% of the total number of clusters. The correction of bias is therefore an important issue in the establishment of investigation priorities, to the extent that it is based on the estimated impact of the clusters.

The remaining gap in the prevention of false positives concerns the beginning (especially December) of the winter viral gastroenteritis epidemic, spread principally by contact at local scale. By their construction, methods A and B each control for season (since the references come from the same season as the target period), but this control is limited because local clusters occurring during the epidemic front can produce spatiotemporal patterns similar to those of waterborne epidemics.

### CONCLUSION

At the conclusion of the pilot phase, it appears that automated searching for waterborne cAGE case clusters is

feasible, and may detect more than 1,000 clusters of five cAGE cases (a mean of 15 AGE cases) or more in France each year, compared to two epidemics of approximately 100 cases of AGE (30 cAGE cases) currently detected by reports of local operators.

In the method presented, the selection of cAGE case clusters is based on the crossing of two complementary methods. This procedure protects against false positives and must be kept for the future. The residual defects in specificity concern the beginning of the winter epidemic period and the possibility of confusion with other types of clusters not otherwise identified (clusters linked to food or swimming).

Other methods of detection, like Kulldorff's space-time scan (Kulldorff *et al.* 2005), will be tested before the national process of automated detection of cAGE case clusters is implemented. The development and evaluation of these new methods of detection is currently underway at InVS. The assessment should continue in the field to determine the capacity to identify foodborne outbreaks and clusters not from waterborne origin, and the feasibility of surveillance of risk factors related to the water supply (Smith *et al.* 2006).

## REFERENCES

- Beaudeau, P. 2010 Natural and human factors of faecal contamination events of drinking water in small distribution networks, France, 2003–2004: a geographical ecological study. *J. Water Health* **8**, 20–34.
- Beaudeau, P., De Valk, H., Vaillant, V., Mannschott, C., Tillier, C., Mouly, D. & Ledrans, M. 2008 Lessons learned from ten investigations of waterborne gastroenteritis outbreaks, France, 1998–2006. *J. Water Health* **6**, 491–503.
- Bounoure, F., Beaudeau, P., Mouly, D., Skiba, M. & Lahiani-Skiba, M. 2011 Syndromic surveillance of acute gastroenteritis based on drug consumption, France. *Epidemiol. Infect.* **139**, 1388–1395.
- Brunkard, J. M., Ailes, E., Roberts, V. A., Hill, V., Hilborn, E. D., Craun, G. F., Rajasingham, A., Kahler, A., Garrison, L., Hicks, L., Carpenter, J., Wade, T. J., Beach, M. J. & Yoder, J. S. 2011 Surveillance for waterborne disease outbreaks associated with drinking water, United States, 2007–2008. *MMWR Surveill. Summ.* **60**, 38–68.
- Delmas, G., Gallay, A., Espie, E., Haeghebaert, S., Pihier, N., Weill, F. X., De Valk, H., Vaillant, V. & Désenclos, J. C. 2006 Les toxi-infections alimentaires collectives en France entre 1996 et 2005 [Foodborne diseases outbreaks in France between 1996 and 2005]. *BEH* **51–52**, 418–422 (in French).
- Frost, F. J., Craun, G. F. & Calderon, R. L. 1996 Waterborne disease surveillance. *Journal AWWA* **88**, 66–75.
- Kulldorff, M., Heffernan, R., Hartman, J., Assuncao, R. & Mostashari, F. 2005 A space-time permutation scan statistic for disease outbreak detection. *PLoS Med.* **2** (3), e59.
- Majowicz, S. E., Hall, G., Scallan, E., Adak, G. K., Gauci, C., Jones, T. F., O'Brien, S., Henao, O. & Sockett, P. N. 2008 A common, symptom-based case definition for gastroenteritis. *Epidemiol. Infect.* **136**, 886–894.
- Rambaud, L., Mouly, D., Schmitt, M., Kerrien, F. & Beaudeau, P. 2011 Utilisation des données de l'Assurance maladie pour évaluer l'impact sanitaire d'une épidémie de gastro-entérites d'origine hydrique, Bourg Saint-Maurice (Arc 1800), 2006 [Use of drug reimbursements data from French national health insurance to characterize a waterborne outbreak in Bourg Saint-Maurice (Arc 1800), France, 2006]. *BEH* **31**, 339–343 (in French).
- Smith, A., Reacher, M., Smerdon, W., Adak, G. K., Nichols, G. & Chalmers, R. M. 2006 Outbreaks of waterborne infectious intestinal disease in England and Wales, 1992–2003. *Epidemiol. Infect.* **134**, 1141–1149.
- Tam, C. C., Rodrigues, L. C. & O'Brien, S. J. 2003 The study of infectious intestinal disease in England: what risk factors for presentation to general practice tell us about potential for selection bias in case-control studies of reported cases of diarrhoea. *Int. J. Epidemiol.* **32**, 99–105.
- Tillet, H. E., De Louvois, J. & Wall, P. G. 1998 Surveillance of outbreaks of waterborne infectious disease: categorizing levels of evidence. *Epidemiol. Infect.* **120**, 37–42.
- Tuppin, P., De, R. L., Weill, A., Ricordeau, P. & Merliere, Y. 2010 French national health insurance information system and the permanent beneficiaries sample. *Rev. Epidemiol. Sante Publique* **58**, 286–290.
- WHO 2005 *Water Safety Plans, Managing drinking water quality from catchment to consumer*. World Health Organization, Geneva.

First received 29 April 2015; accepted in revised form 25 October 2015. Available online 4 December 2015