**Additional file2 Summary of cleaning process and data set exploration (data set 1)**

| Data set 1 | Observations *(%)* | Action |
|---|---|---|
| **A**      *Cleaning process 1* | *N = 3187* | |
|      *Select the study period* | | |
| Year 2012 | 48 *(1.5)* | removed |
| **B**      *Cleaning process 2* | *N = 3139* | *saved* |
|      *Check multiple entries* | | |
| "Equivocal" and "negative" *Test Results* | 532 *(16.9)* | removed |
| Double entry of individuals | 30 *(0.9)* | removed |
| **C**      *Cleaning process 3* | *N = 2576* | *saved* |
|      *Information about the data set* | | |
| Epi-linked cases | 60 *(1.9)* | kept |
| Missing *Age* | 6 *(0.2)* | kept |
| Missing *Onset Date* | 0 *(-)* | - |
| "Indeterminate", "not stated/unknown" *Gender* | 18 *(0.7)* | kept |
| Missing *Serogroup code* | 483 *(18.8)* | kept |
| *Country of Source* "not stated" and "Inadequately described" | 481 *(18.7)* | kept |
| Missing *locality name* | 0 *(-)* | - |
| Misspelling *Street Name* | 107 *(4.1)* | corrected |
| *Street Number* in the wrong cell | 5 *(0.2)* | corrected |
| Missing *Street Name* | 258 *(10.0)* | kept |
| Missing *Street Number* | 309 *(11.9)* | kept |
| **D**      *Cleaning process 4* | *N = 2576* | *saved* |
|      *Preparation for Time series data set for locally –acquired cases* | | *Table 2* |
| Locally-acquired cases | 2418 *(93.9)* | kept |
| CCD code not attributed | 306 *(11.9)* | kept |
|     SLA name not stated or 'Overseas - other" | 132 *(5.1)* | removed |
|     "Indeterminate" *Age* | 5 *(0.2)* | removed |
|     "Indeterminate" *Gender* | 19 *(0.7)* | removed |
| **E**      *Cleaning process 5* | *N= 2262* | |
|      *Develop the time-series dataset by :* | | |
| SLA | 43 | |
| Age-group | 17 | |
| Gender | 2 | |
| Year | 17 | |
| Month | 12 | |
| **Total time-series data set** | *= 298248* | saved and analysed |