



A time-series approach for clustering farms based on slaughterhouse health  
aberration data

Hulsegge, I., de Greef, K. H., & Hulsegge, I.

This is a "Post-Print" accepted manuscript, which has been published in "Preventive  
Veterinary Medicine"

This version is distributed under a non-commercial no derivatives Creative Commons



([CC-BY-NC-ND](https://creativecommons.org/licenses/by-nc-nd/4.0/)) user license, which permits use, distribution, and reproduction in any medium, provided the original work is properly cited and not used for commercial purposes. Further, the restriction applies that if you remix, transform, or build upon the material, you may not distribute the modified material.

Please cite this publication as follows:

Hulsegge, I., de Greef, K. H., & Hulsegge, I. (2018). A time-series approach for clustering farms based on slaughterhouse health aberration data. *Preventive Veterinary Medicine*, 153, 64-70. DOI: 10.1016/j.prevetmed.2018.03.003

You can download the published version at:

<https://doi.org/10.1016/j.prevetmed.2018.03.003>

1 **A time-series approach for clustering farms based on**  
2 **slaughterhouse health aberration data**

3 B. Hulsegge<sup>a</sup> and K.H. de Greef<sup>a</sup>

4

5 <sup>a</sup> Animal Breeding and Genomics, Wageningen Livestock Research, P.O. Box 338, 6700 AH,  
6 Wageningen, The Netherlands.

7

8 Corresponding Author:

9 Ina Hulsegge, Animal Breeding and Genomics, Wageningen Livestock

10 Research, P.O. Box 338, 6700 AH, Wageningen, The Netherlands. Tel: +31-317-480513; E-

11 mail: [ina.hulsegge@wur.nl](mailto:ina.hulsegge@wur.nl)

12

13 **Abstract**

14 A large amount of data is collected routinely in meat inspection in pig slaughterhouses. A  
15 time series clustering approach is presented and applied that groups farms based on  
16 similar statistical characteristics of meat inspection data over time. A three step  
17 characteristic-based clustering approach was used from the idea that the data contain  
18 more info than the incidence figures. A stratified subset containing 511,645 pigs was  
19 derived as a study set from 3.5 years of meat inspection data. The monthly averages of  
20 incidence of pleuritis and of pneumonia of 44 Dutch farms (delivering 5,149 batches to 2  
21 pig slaughterhouses) were subjected to 1) derivation of farm level data characteristics 2)  
22 factor analysis and 3) clustering into groups of farms. The characteristic-based clustering  
23 was able to cluster farms for both lung aberrations. Three groups of data characteristics  
24 were informative, describing incidence, time pattern and degree of autocorrelation. The  
25 consistency of clustering similar farms was confirmed by repetition of the analysis in a  
26 larger dataset. The robustness of the clustering was tested on a substantially extended  
27 dataset. This confirmed the earlier results, three data distribution aspects make up the  
28 majority of distinction between groups of farms and in these groups (clusters) the  
29 majority of the farms was allocated comparable to the earlier allocation (75% and 62%  
30 for pleuritis and pneumonia, respectively). The difference between pleuritis and  
31 pneumonia in their seasonal dependency was confirmed, supporting the biological  
32 relevance of the clustering. Comparison of the identified clusters of statistically  
33 comparable farms can be used to detect farm level risk factors causing the health  
34 aberrations beyond comparison on disease incidence and trend alone.

35

36 **Highlight**

- 37
- Characteristic-based clustering is able to cluster time series of meat inspection  
38 data of farms using a set of derived statistical characteristics.
  - Seasonality and data dispersion characteristics such as autocorrelation have  
39 additional value to the conventional incidence figures of pneumonia and pleuritis.  
40

- 41       • Farms were mainly clustered on: amount of variation in the data; distribution  
42       shape of the data and similarity between consecutive data points.

43

#### 44 **Keywords**

45 Meat inspection data; Time series; Characteristic-based clustering; Big data ; Pneumonia  
46 ; Pleuritis

47

#### 48 **1. Introduction**

49 According to legal regulations (European Community, 2004), all slaughtered pigs in the  
50 European Union are subject to a routine meat inspection at the slaughterhouses.

51 Traditionally, meat inspection has been used to reduce food-borne risk to public health  
52 (Edwards et al., 1997). The meat inspection findings are also valuable indicators that can  
53 be used as a feedback system indicating animal health and to derive recommendations  
54 for improvement of farm management (Schuh et al., 2000). Meat inspection data can be  
55 used to inform farmers on the health status of their herd (benchmarking) since health  
56 aberrations indicate systems (housing, ventilation control) or management (treatment  
57 and prevention strategies) failures. Slaughterhouse data both reveal such problems and  
58 offer the opportunity to monitor effectivity of interventions. Current use of  
59 slaughterhouse health aberration data seems limited to periodic reporting of farm  
60 incidence averages. Understanding the data structure (such as temporal patterns) of  
61 aberrations in meat inspection data may provide important information beyond these  
62 average incidence figures.

63 One possible approach to analyse meat inspection data involves time series methods,  
64 such as exploratory methods (Sanchez-Vazquez et al., 2012; Alhaji et al., 2015) and  
65 autoregressive models (Neumann et al., 2014; Vial and Reist, 2014; Adachi and Makita,  
66 2015). These methods however, require structured data, a sufficient number of  
67 observations that are fairly regularly measured over time, which is often not the case for  
68 data on batches of pigs delivered to slaughterhouses. Another possible approach is time

69 series clustering directly on raw data. This method however does not account for the  
70 temporal sequences of the observed values and the autocorrelations structure of the data  
71 is ignored. Characteristic-based clustering has been developed to address the problem of  
72 clustering raw time series data (Hennig et al., 2015). This method has been proposed by  
73 several authors in various domains such as electricity (Räsänen and Kolehmainen, 2009),  
74 business (Davenport and Funk, 2015), and human health (Leffondré et al., 2004)  
75 (Niedermeyer et al., 2011). We applied this method to group farms based on similar  
76 statistical characteristics of meat inspection data, focussing on pneumonia and pleuritis.  
77 The objective of this study was to explore whether an analysis which utilises more  
78 information from the data than incidence figures provides added value to make  
79 distinctions between individual farms. A comprehensive meat inspection dataset,  
80 collected over 3.5 years, was available for this. This more detailed farm characterisation  
81 may aid in finding risk factors for failures by comparing more uniform groups of farms.

82

## 83 **2. Material and methods**

### 84 2.1. Data Source

85 Post mortem meat inspection data of carcass and organs are collected on every  
86 slaughtered pig in The Netherlands. The inspection procedures are described in detail in  
87 Regulation EC no. 854/2004 (European Community, 2004). Meat inspection data  
88 collected between January 2011 and August 2014 were provided by the major Dutch  
89 meat producer, one record for each slaughtered pig, with information on pneumonia and  
90 pleuritis and aberrations on legs, skin and liver. Respiratory disorders were chosen as  
91 study dataset as they are one of the major diseases affecting pigs worldwide (Brockmeier  
92 et al., 2002) and have reasonable incidences across farms and seasons and the  
93 repeatability of the slaughterhouse classification is adequate.

94

### 95 2.2. Study sample

96 Criteria were developed to derive a suitable sub-dataset for method development and  
97 analysis. August 2014 was excluded since it did not comprise the entire month, also

98 batches with less than 10 animals were excluded. The two slaughterhouses with the  
99 largest number of records were selected. These slaughterhouses had complete datasets  
100 for the entire period and no obvious changes in inspection system. In this set, farms  
101 were selected that had delivered at least one batch with at least 10 pigs every month and  
102 at least 87 batches (more than 1 batch per 2 weeks on average). The resulting study  
103 sample contained information of 511,645 pigs submitted from 44 Dutch farms in 5,149  
104 batches.

105 Information on the percentage pneumonia and pleuritis in the batches is presented in  
106 Table 1. The analysis is principally batch based – records were created containing batch  
107 averages. The percentage of each aberration (pleuritis or pneumonia) in each batch was  
108 computed as number of pigs in that batch with the aberrations divided by total number of  
109 pigs in that batch multiplied by 100.

110

111 The study dataset is quite complete from a statistical point of view (no missing records,  
112 good distribution over the entire study period), but comprises a small part of the total  
113 dataset. For verification and validation reasons a second, larger, dataset was created.  
114 The selection criteria were released: all farms of the two slaughterhouses were included  
115 which met the criterion that the whole study period (all months) was reasonably covered:  
116 6 month averages were allowed to be missing for each farm. This resulted in an three to  
117 almost fourfold size of the data: 163 farms delivering 15,276 batches comprising  
118 1,829,762 slaughtered pigs. Table 1 illustrates that the characteristics of the validation  
119 set resemble those of the study set.

120

121

122 Table 1. Percentage pneumonia and pleuritis in the study sample (5,149 batches) and  
 123 validation sample (15,276 batches).

Aberration	# Batches with percentage 0% (%)	Mean percentage (95% CI) in a batch	Sd percentage	Max percentage
Pneumonia				
Study sample	615 (11.9%)	8.76 (8.51–9.01)%	9.10%	63.83%
Validation sample	1599 (10.5%)	9.12 (8.96 – 9.26)%	9.38%	78.15%
Pleuritis				
Study sample	375 (7.3%)	12.42 (12.12-12.72)%	10.91%	61.64%
Validation sample	1384 (9.06%)	10.04 (9.88-10.20)%	10.37%	82.58%

124

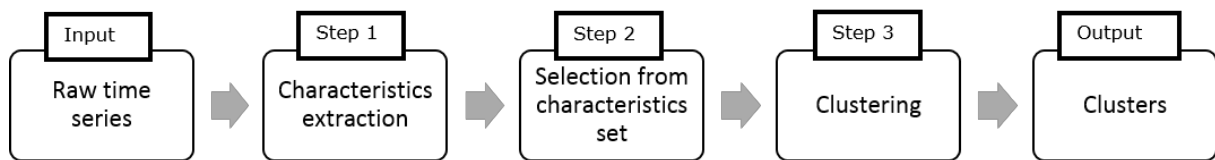
125       2.3. Time series visual explorations

126 For exploratory purpose, percentage aberrations were aggregated for each month of  
 127 study. An exploratory analysis was conducted by plotting percentage aberrations of the  
 128 study sample containing 44 farms in the period January 2011 to July 2014 in a  
 129 multivariate time series plot using the R package mvtsplot (Peng, 2008). The mvtsplot  
 130 method produces an adaptation of the multivariate time series plot which combines a  
 131 heatmap with boxplot-like summaries and a basic line plot to provide a detailed overview  
 132 of the data. The colours purple, grey and green in the heatmap correspond to low,  
 133 medium and high values, respectively. The darker the shading the larger the value.

134

135       2.4. Time series clustering using global characteristics

136 We used a three step method to group farms with comparable statistical characteristics  
 137 of health aberrations over time (Fig. 1). The first step of the method involved replacing  
 138 the raw time series data with some global measures of time series characteristics, as  
 139 described by Wang et al. (2006) and Räsänen and Kolehmainen (2009). The measures  
 140 summarized information of the time series, to capture the ‘global picture’ of the data.



141 Fig. 1. Characteristics based clustering approach (after Wang et al. (2006)).  
142

143  
144 The characteristics used in this study were: *mean, standard deviation, trend, seasonality,*  
145 *remainder, autocorrelation, skewness, kurtosis, chaos, nonlinearity, and self-similarity.*

146 Table 2 describes the popularised interpretation of these characteristics and their  
147 acronym used below.

148 *Trend* and *seasonality* are common characteristics of time series, and it is natural to  
149 characterize a time series by its degree of trend and seasonality. In addition, once the  
150 trend and seasonality of a time series has been measured, the time series can be  
151 detrended and deseasonalised to enable additional features such as noise or chaos to be  
152 more easily detectable. The R function *stl* was used for detrending and deseasoning the  
153 timeseries (Cleveland et al., 1990). For the validation sample (which contained missing  
154 values), the R package *stlplus* version 0.5.1 was used to detrend and deseasonalise the  
155 time series, applying a loess algorithm to handle missing values (Hafen, 2010).

156 To obtain a precise and comprehensive calibration, some measures are calculated on  
157 both the raw time series as well as the remaining time series after detrending and  
158 deseasonalising. All these characteristics (presented in a popular phrasing in table 2) are  
159 thoroughly explained by Wang et al. (2009) and (Davenport and Funk, 2015).

160

161



162 Table 2. Summary of the used data characteristics, calculated from the raw batch data  
 163 and on the detrended and deseasonalised data.

Characteristic	Definition	Acronym Raw data	Acronym Detrended and deseasonalised data
Mean	The average of the observations	' <i>mean</i> '	
Standard deviation	A measure of how spread out the data is.	' <i>sd</i> '	
Trend	A pattern found in time series; used to describe whether the data is showing an upward or downward movement for a part, or all of the time series.	' <i>trend</i> '	
Seasonality	A pattern of a time series in which the data experience regular and predictable changes that repeat every calendar year.	' <i>seasonality</i> '	
Remainder	The residuals of the time series after allocation into the seasonal and trends time series (also called "noise", "irregular" or "random").	' <i>remainder</i> '	
Hurst Exponent	A measure for longterm memory and fractality of a time series (an evaluation index of the self-similarity).	' <i>self.sim</i> '	
Autocorrelation	The correlation within a time series with its own past and future values (also called serial correlation)	' <i>autocorr</i> '	' <i>dc-autocorr</i> '
Skewness	A measure of how symmetrical a distribution is.	' <i>skewness</i> '	' <i>dc-skewness</i> '
Kurtosis	A measure which describes the distribution of the observed data around the mean. A measure of how peaked or flat a distribution is relative to the normal distribution	' <i>kurtosis</i> '	' <i>dc-kurtosis</i> '
Lyapunov Exponent	A measure of stability; Chaos	' <i>chaos</i> '	
Nonlinearity	A measure for not arranged in a straight line.	' <i>nonlin</i> '	' <i>dc-nonlin</i> '

164

165 In the second step a factor analysis, using the function *principal* from the R package

166 *Pysch* version 1.5.8 (Revelle, 2015), was performed to select a subset of characteristics

167 that condensed the information present in the characteristics and provided the best  
168 description. We only kept the factors with an eigen-value greater than 1 (Tabachnick and  
169 Fidell, 2006), those that are more informative than a single variable. The varimax  
170 rotation was used to facilitate the interpretation of results by maximising the loading of  
171 each individual variable on a single factor (i.e., its correlation with this factor). For each  
172 factor the measure that had the highest loadings (i.e. the highest correlation with a give  
173 factor) was selected.

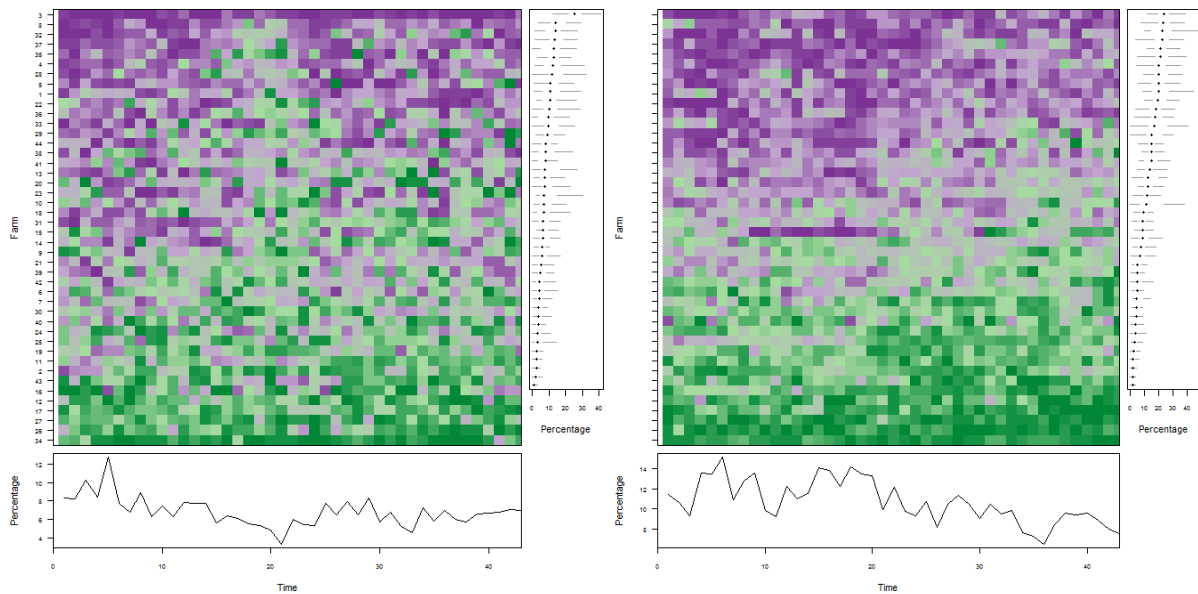
174 Finally (third step), we used cluster analysis to identify clusters of farms with similar  
175 patterns of characteristics selected by the factor analysis. In order to weigh all  
176 characteristics equally, all characteristics were transformed to the same range (0,1). A  
177 measure near 0 for a certain time series indicates an absence of the characteristic while a  
178 measure near 1 indicates a strong presence of the characteristic (Wang et al., 2006). The  
179 measures were normalised with the function *SofMax* of the R package *DMwR* version  
180 0.4.1 (Torgo, 2010). The R package *NbClust* version 3.0 (Charrad et al., 2014) was used  
181 to perform the cluster analysis, in order to identify the optimal number of clusters.  
182 Clusters were generated using the complete linkage method applied to Euclidean  
183 distances.

184  
185

### 186 **3. Results**

#### 187 3.1. Percentage aberrations at farm level.

188 The monthly percentage of aberrations for the farms in the periods January 2011 to July  
189 2014 varied between farms and months (Fig. 2). Monthly farm pneumonia incidences in  
190 the study set varied between 0.8 and 25.1%, and pleuritis incidences varied between 1.4  
191 - 24.0%. As the colouring in Figure 2 indicates, farms not only differ in monthly  
192 incidence, there is also considerable between farm variation in within-farm consistency in  
193 time. Consistent farms have either an entirely green coloured time-series (low incidence)  
194 or entirely purple coloured time-series (high incidence). Farms with alternating colours in  
195 their time series have low consistency in their incidences.



196

197 Fig. 2. Multivariate time series plot of percentage pneumonia and pleuritis for 44 farms.

198 The purple to green palette represents variation in percentage aberrations (green

199 represents low percentages; purple high percentages). The right panel presents

200 summary statistics of percentage aberrations for each farm, the black dots denote the

201 median while the horizontal lines represent the lower and upper quartiles. The lower

202 panel shows the median values of percentage aberration across the time series of the 43

203 months (1= January 2011 and 43= July 2014) for each time point.

204

205

206 3.2. Time series clustering using global characteristics.

207 3.2.1. Pneumonia

208 Exploratory factor analysis of percentage pneumonia reduced the 15 global

209 characteristics to four factors explaining 65% of the variance. The most informative

210 global characteristics were: '*mean*', '*seasonality*', '*autocorr*' and '*dc-kurtosis*' (Fig. 3).

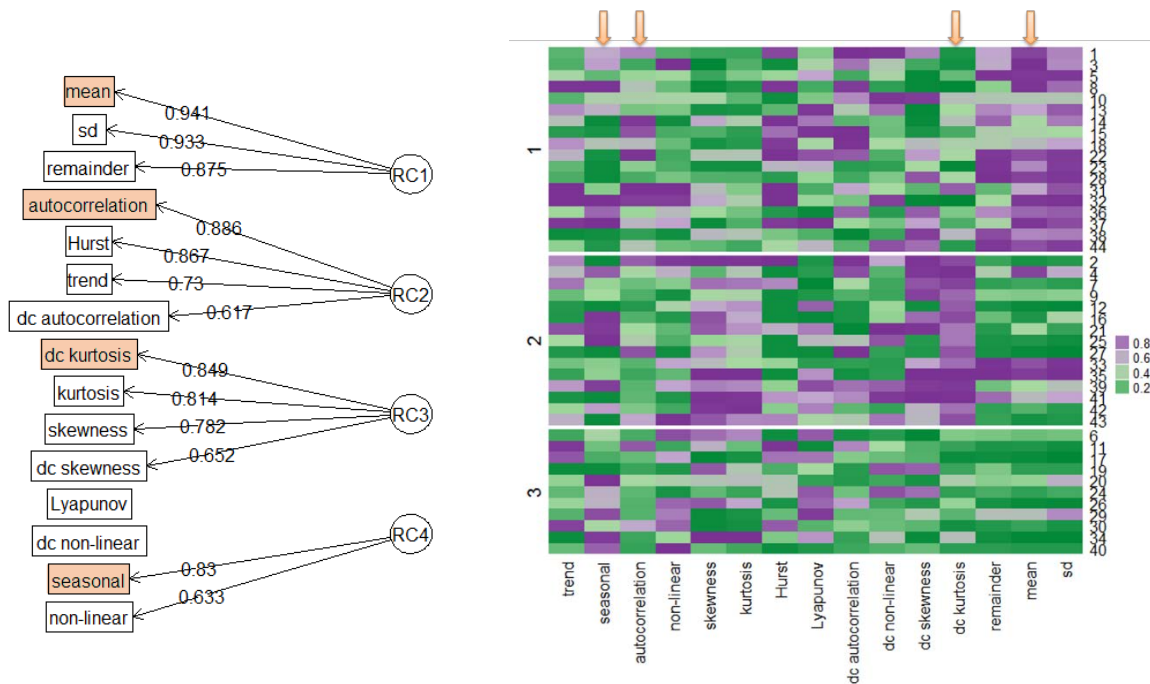
211 These four most informative global characteristics were used by cluster analysis and

212 resulted in grouping the 44 farms into three clusters. Categorization of the three clusters

213 data characteristics are shown in Figure 3.

214

215



216

217 Fig. 3. Factor analysis path diagram of pneumonia (left pane) and allocation to three  
 218 clusters (right pane). Left pane: The coloured square boxes are the characteristics of  
 219 each principal component (PC) that are used in subsequent analysis. On the straight  
 220 arrows, the loadings (correlation between the principal component and the characteristic)  
 221 are presented. Only the largest loadings are shown. Right pane: Characteristics summary  
 222 for the three identified clusters. Purple indicates high values, green indicates low values.  
 223 The arrows at the top indicate the selected global characteristics; left axis: cluster  
 224 number; right axis: farm number.

225

226 Farms in cluster 1 are characterised by high incidence values (*'mean'*) with large  
 227 variability in pneumonia incidence and low trend and seasonally adjusted kurtosis (*'dc-*  
 228 *kurtosis'*), having a flat top near the mean and produces fewer and less extreme outliers  
 229 than does the normal distribution. Cluster 2 groups farms with the opposite: low *'mean'*  
 230 and high *'dc-kurtosis'*, having a distinct tall peak near the mean, decline rather rapidly  
 231 and have fatter tails or more extreme values. Farms in cluster 3 share the low incidences  
 232 with cluster 2, but combine this with low kurtosis, meaning that the trend and seasonally  
 233 adjusted time series produces fewer and less extreme outliers than does a normal  
 234 distribution. The factor analysis suggested *'seasonality'* as an informative characteristics,

235 the value for all clusters showed little recurring seasonal pattern, periods of above-  
236 average and below-average percentage pneumonia each year (Fig. 3). Farms belonging  
237 to cluster 1 fluctuated most with season, from -2.4% in September to 2.2% in  
238 December. For cluster 2 the lowest value of the seasonal component was observed in  
239 August (-1.1%) and the highest in May (1.2%). For cluster 3 these values varied from -  
240 1.4% in September to 1.5% in May.

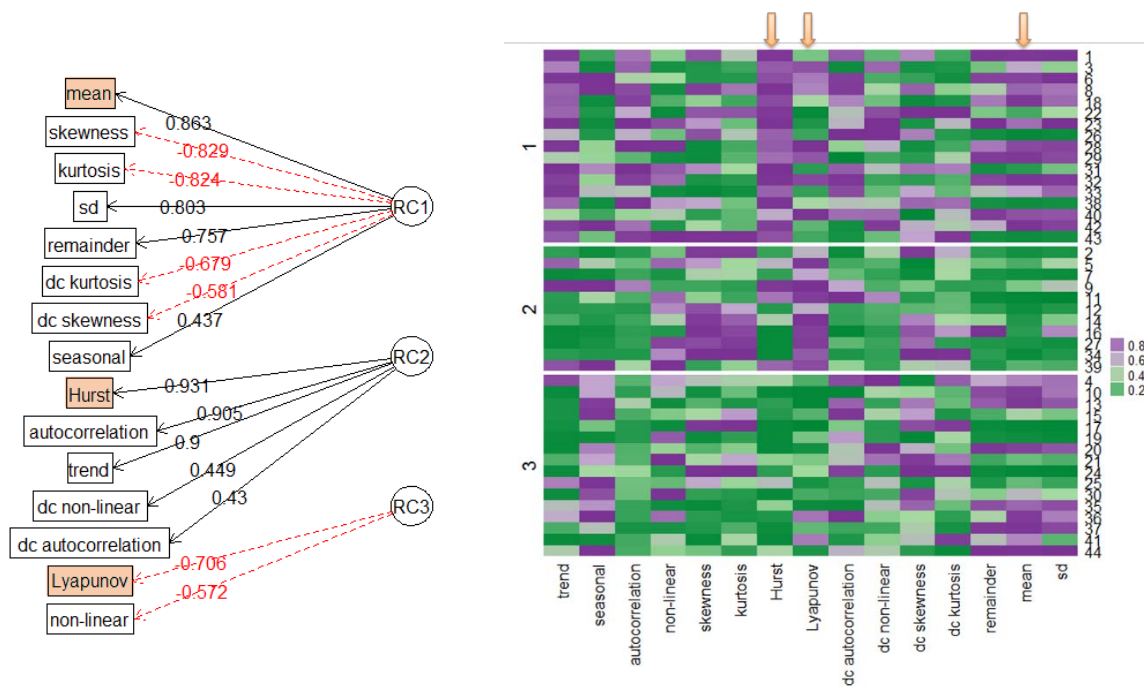
241

### 242 3.2.2. Pleuritis

243 Factor analysis of percentage pleuritis reduced the 15 global characteristics to three  
244 factors with the most informative global characteristics being '*self.sim*', '*chaos*' and  
245 '*mean*', explaining 62% of the variance (Fig. 4). Cluster 1 contains farms with high  
246 incidence values and a strong trend. These time series had also highly regular  
247 fluctuations over time (high '*self.sim*' exponent values; indicating a persistent time  
248 series) and showed no chaotic behaviour (low '*chaos*' values). Farms with low incidence  
249 figures were predominantly allocated to cluster 2, combined with high values for '*chaos*'.  
250 The time series showed no trend or seasonal effect and self-similarity was almost not  
251 present. Farms in cluster 3 are characterised by low levels of '*autocorr*', '*chaos*' and  
252 '*self.sim*', but differ mutually in their incidence figures.

253

254



255  
 256 Fig. 4. Factor analysis path diagram of pleuritis (left pane) and allocation to three  
 257 clusters (right pane). For explanation: see Figure 3.

258  
 259 3.2.3. Validation in the extended dataset

260 To test the robustness of the clustering, the analysis was repeated with a substantially  
 261 extended dataset. By releasing 1) the criterion that farms have to have batches in all  
 262 months of the study period and 2) the criterion that at least 87 batches were delivered,  
 263 the study size was extended about fourfold. The only remaining criteria were that  
 264 batches contained at least 10 pigs and that in most months a batch was available,  
 265 maximally six missing. The *stlplus* allows handling of missing monthly values.

266  
 267 For both pneumonia and pleuritis aberration data, again the three archetypes of clusters  
 268 evolve: one based on the variation characteristics, one on the distribution shape and one  
 269 on the similarity of consecutive data points (Fig. A1 and A2). Specifically for pneumonia,  
 270 the validation set, having missing data points in most farms, failed to identify the specific  
 271 cluster indicating a group of farms with specific seasonal sensitivity. Similarly, the  
 272 validation exercise on pleuritis figures was to some degree less distinctive in  
 273 discriminating farms with regard to sequentiality/chaos, but was stronger in its

274 separation between variance (*'mean', 'sd'*) and shape characteristics (*'skewness',*  
 275 *'kurtosis'*) with regard to pleuritis.

276

277 In Table 3 and 4, the degree of similarity in allocation to the clusters between the study  
 278 analysis and the validation analysis is presented for the 44 farms that were involved in  
 279 the study sample and again allocated to new clusters during from the validation analysis.

280

281 Table 3. Clustering of farms for pneumonia characteristics: comparison of the coherence  
 282 in allocation to clusters in the test analysis (vertical: clusters 1, 2, 3) and to clusters in  
 283 the validation analysis (horizontal: clusters A, B, C, D).

Validation cluster →	A	B	C	D	Total
Original cluster ↓					
1	3	15			18
2	3	4	1	7	15
3	3	5		3	11
Total	9	24	1	10	44

284

285

286 Table 4. Coherence in clustering of farms for pleuritis characteristics: number of farms  
 287 allocated to clusters in the test analysis (vertical: clusters 1, 2, 3) and to clusters in the  
 288 validation analysis (horizontal: clusters A, B, C).

Validation cluster →	A	B	C	Total
Original cluster ↓				
1	4	13		17
2	4	1	6	11
3	14		2	16
Total	22	14	8	44

289

290 For pleuritis, the overall correspondence of the two clustering analyses amounts

291 75%. 14, 6 and 13 farms (from cluster 1, 2 and 3 respectively) which were grouped  
292 together in the original analysis were again allocated together into the new clusters. For  
293 pneumonia, the overall correspondence is less. It amounts 61% and an extra cluster is  
294 formed. The largest new cluster (B, containing 24 farms) comprises the majority of two  
295 of the original clusters (1 and 3). Overall, most of the farms that were grouped together  
296 in the study analyses were allocated into joint clusters again in the validation study, both  
297 for pleuritis and pneumonia.

298

#### 299 **4. Discussion**

300 Meat inspection generates a large amount of time series data that are used to only a  
301 limited extend for animal health surveillance purposes. And if so, use is generally limited  
302 to the average incidence (*'mean'*), and its change in time, solely on farm level. Current  
303 exercise enriches this by combining data across farms. Understanding the underlying  
304 information and interpretation of the results for meaningful purposes (such as  
305 management support or detection of risk factors) is an opportunity, but also a challenge  
306 due to the high diversity between farms, batches and underlying factors.

307 A dataset containing more than 3½ years of historical meat inspection data was available  
308 to explore the potential of a data analysis to cluster farms into groups with comparable  
309 health aberration patterns over time. In this dataset two respiratory disorders,  
310 pneumonia and pleuritis, were chosen as study objects as they are among the major  
311 diseases affecting pigs worldwide (Brockmeier et al., 2002), (Merialdi et al., 2012) (Eze  
312 et al., 2015) and the most common slaughter aberrations found in pigs (Sanchez-  
313 Vazquez, 2013). Also, the diversity of incidences and the reasonable repeatability of  
314 slaughterhouse pleuritis and pneumonia classification is helpful from a statistical point of  
315 view to develop the proposed method.

316

##### 317 **4.1. The method**

318 Clustering is among the most widely used method in the analysis of time series data  
319 (Fidaner et al., 2015) (Chen et al., 2017) and for our casus, it offers the opportunity to



320 identify farms with similar patterns of percentage pneumonia or pleuritis over time  
321 discerning similarities between those farms beyond obvious characteristics such as  
322 incidence figures (Fidaner et al., 2015).

323 Characteristic-based clustering first converts raw time series data into a characteristic  
324 vector of lower dimension, after which clustering is applied. Characteristic-based  
325 clustering, in the literature also called Feature based clustering or Statistical measures  
326 based clustering, has been proposed by several authors across science for clustering time  
327 series. For example, Leffondré et al.(2004) used this method for identifying patters of  
328 change in quantitative human health indicators and Räsänen and Kolehmainen (2009) for  
329 electricity use time series data. We applied this method to group farms based on similar  
330 statistical characteristics of meat inspection data over time. Our approach consists of  
331 three distinct steps: 1) computation of the data characteristics on farm level from  
332 monthly farm averages; 2) factor analysis to identify the major explaining variation  
333 among farms; and 3) cluster analysis to group farms on basis of similarity in their data  
334 characteristics. Ad 1), we used the set of characteristics as proposed by Wang et al.  
335 (2005; 2006) that contains measures of '*trend*', '*seasonality*', '*autocorr*', '*skewness*',  
336 '*kurtosis*', '*chaos*', '*nonlin*', and '*self.sim*' to represent time series. The proposed  
337 statistical characteristics were selected because they are simple and easy to compute. Ad  
338 2), for selection of the most relevant characteristics of the data set, various approaches  
339 can be used. We used factor analysis as search mechanism to find the best selection  
340 from the characteristics set as suggested by Leffondré et al. (2004). This is an easy and  
341 widely accepted method to identify common patterns in data with diverse correlations  
342 structures. Ad 3), we chose the clustering method according to Leffondré et al. (2004) as  
343 it seems to fit our ambition well.

344

345 The correlation matrix between the characteristics illustrates that several characteristics  
346 were highly correlated; e.g. mean and standard deviation for percentage pneumonia as  
347 well as percentage pleuritis (data Fig. A3 and A4). From a methodological point of view,  
348 this correlation structure implies that some features are interchangeable. Having two

349 highly correlated characteristics makes one virtually redundant – in this case it may be  
350 useful to select the one which is easiest to interpret, as suggested by Leffondré et al.  
351 (2004).

352

#### 353 4.2. The clustering of farms on basis of statistical characteristics

354 The results showed that the applied approach is able to discriminate between farms with  
355 regard to their meat inspection data. Both the data on pneumonia and on pleuritis  
356 resulted in three clusters. A closer look into the composition of these clusters reveals that  
357 both in the pleuritis data and in the pneumonia data farms were clustered mainly on: 1)  
358 amount of variation in the data; 2) distribution shape of the data and 3) similarity  
359 between consecutive data points. Both the consistent distinction between groups of  
360 characteristics and the consistency between the study results and the validation results  
361 confirm that the method is able to make distinction between farms beyond grouping  
362 them on the conventional way: incidence (percentage of pigs), possibly grouped in  
363 categories like high, moderate and low incidence.

364

#### 365 4.3. Study set versus validation set.

366 The dataset on 44 farms (511,645 animals) was the ideal set to develop the method.  
367 But, regarding the small sample size, quite distant from the data as a whole. Extension  
368 to a larger (163 farms, 1,829,762 animals), but less optimal (less data points per farm,  
369 some missing month averages) set is a feasible model to verify whether the method  
370 holds for in a less ideal situation. This confirmed the earlier results, three data  
371 distribution aspects make up the majority of distinction between groups of farms and in  
372 these groups (clusters) the majority of the farms was allocated comparable to the earlier  
373 allocation (75% and 62% for pleuritis and pneumonia, respectively).

374

375 Switching to less structured data also revealed a trade off between accuracy (a small but  
376 precise data set) and volume (a larger but more rough dataset). The study sample

377 revealed a specific vector for seasonal sensitivity for pneumonia, which was not detected  
378 in the larger dataset which had missing datapoints.

379

#### 380 4.4. Relevance

381 Classically, farms are compared on basis of the incidence of lung problems. Obvious first  
382 next level comparisons comprise the variability and change in time of individual farm  
383 health performance. The high correlation between mean, standard deviation and  
384 coefficient of variation of the aberrations in both pneumonia and pleuritis indicates that  
385 variability between batches is not a valuable extra trait in itself, as it does not add  
386 substantial information additional to the average level of aberrations. On the other hand,  
387 other characteristics, such as repeatability patterns in time do aid in making distinctions  
388 between farms.

389

390 In literature, farm factors that affect problems like pleuritis and pneumonia are often  
391 assessed by comparing farming systems factors such as organic versus conventional  
392 farms (e.g. Alban et al. (2015) or comparing large versus small scale farms (e.g. (Enoe  
393 et al., 2002; Fablet et al., 2012). Data analysis offers an additional entry: the statistical  
394 grouping of farms may point at similarities in farm characteristics within the groups or  
395 differences between the groups (clusters) of farms that do not vary between for example  
396 organic and conventional systems, but rather are underlying factors in both systems that  
397 are causally related to the incidence of health aberrations. The clustering approach thus  
398 goes beyond comparing farms on basis of systems characteristics (size, type) or  
399 performance (incidences of aberrations) alone and bears the promise to reveal relevant  
400 risk factors from data of seemingly similar farms

401

402 A real practical validation requires insight of the farm characteristics. Relating farm  
403 characteristics (farm size, housing characteristics etc.) to the clusters is the next step to

404 utilize its relevance for enhancing health performance. A promising approach to identify  
405 risk factors for lung aberrations is to study whether the farms in different clusters also  
406 structurally deviate at farm-level either in (nutritional) management practices or in  
407 environmental (housing and ventilation) factors. Our dataset was unique in its size and  
408 consistency, but it contains only slaughter data, farms were coded, implying that no farm  
409 characteristics were available in the analysis. On availability of adequate data, comparing  
410 the clusters with regard to farm characteristics is an obvious next step in studying added  
411 value of these clustered slaughterhouse data. Do farms that are clustered on statistical  
412 grounds also resemble in farm characteristic? And which farm characteristics? If so, this  
413 is a signal that these characteristics may be closely related to real risk factors. Further  
414 developments in data sharing and in data analytics (*big data, machine learning*) are likely  
415 to further develop such opportunities.

416

#### 417 4.5. Biological interpretation/ relevance

418 Pleuritis and pneumonia are both disorders of the respiratory system, but have different  
419 aetiology. Present paper is not intended to elaborate on this, but lines towards biological  
420 interpretation can be drawn. Patterns and trends in incidences of pleuritis and pneumonia  
421 are readily discernable in massive slaughterhouse data, but are difficult to quantify in  
422 detail on the individual farm level. Current method identifies these patterns on individual  
423 farms, making use of the trends and patterns of related farms. And it subsequently  
424 groups farms with similar aberration characteristics. This grouping is considerably  
425 different for pneumonia and for pleuritis. Comparison of the clusters of farms for  
426 pneumonia and for pleuritis reveals that only a minority of the farms shares the same  
427 clusters (data not shown). This illustrates the different underlying factors affecting the  
428 (slaughterhouse detected pathological indicators of) these two respiratory disorders.  
429 Furthermore, the analysis identified within- and between farm variation related to  
430 seasonality for pneumonia, rather than for pleuritis. This is in line with earlier studies (for  
431 example by Fablet et al. (2012) who identified distinct risk factors for pneumonia and

432 pleuritis (ventilation and seasonality versus temperature and barn climate) in slaughtered  
433 pigs.

434 Interpretation of the parameters from a biological point of view is possible, but  
435 speculative. For example, the parameters for seriality may indicate that farms with high  
436 figures for this have the characteristic that they are quite consistent between months in  
437 their aberration performance. Also, high or low levels of kurtosis could be interpreted as  
438 relatively long or short problem periods. However, such interpretations are speculative,  
439 and to our knowledge, such interpretations have not been made in literature.

440 Results like those presented here confirm the biological ground under the identified  
441 clusters. Also, they support the expectation that data analysis points at less obvious  
442 underlying phenomena, which is helpful in further understanding the farm level aetiology  
443 of these disorders. Also, management opportunities can be strengthened by combing  
444 farm characteristics to the wealth of routinely collected data in slaughter houses,  
445 primarily in detecting husbandry related risk factors.

446 Current work has higher relevance for practical application (such as identification of farm  
447 factors affecting incidence levels) than for enhanced understanding of the underlying  
448 biology of the diseases involved.

449

## 450 **5. Conclusion**

451 Characteristic-based clustering was able to cluster time series of meat inspection data of  
452 farms using a set of derived statistical characteristics. The stepwise analysis of the  
453 slaughterhouse dataset reveals structured variation among farms in incidence of  
454 pneumonia and pleuritis. The applied method groups them into clusters of 'similar' farms  
455 beyond clustering them just on basis of observed incidence of aberrations. Seasonality  
456 and data dispersion characteristics such as autocorrelation had additional value to the  
457 conventional disease incidence figures. The differences between the clusters likely point

458 at systematic differences between individual farms. Validation on a substantially  
459 extended dataset confirmed the results of the study dataset.

460

## 461 **Acknowledgements**

462 This research was commissioned and partly funded by The Feed4Foodure program line  
463 "Nutrition, Intestinal Health, and Immunity" (Feed4Foodure; BO-22.04-002-001). The  
464 authors would like to thank Vion Food Group for providing the data.

465

## 466 **References**

- 467 Adachi, Y., Makita, K., 2015. Real time detection of farm-level swine mycobacteriosis  
468 outbreak using time series modeling of the number of condemned intestines in  
469 abattoirs. *The Journal of Veterinary Medical Science* 77, 1129-1136.
- 470 Alban, L., Petersen, J.V., Busch, M.E., 2015. A comparison between lesions found during  
471 meat inspection of finishing pigs raised under organic/free-range conditions and  
472 conventional, indoor conditions. *Porcine Health Management* 1, 4.
- 473 Alhaji, N.B., Odetokun, I.A., Shittu, A., Onyango, J., Chafe, U.M., Abubakar, M.S.,  
474 Muraina, I.A., Fasina, F.O., Lee, H.S., 2015. Time-series analysis of ruminant  
475 foetal wastage at a slaughterhouse in North Central Nigeria between 2001 and  
476 2012.
- 477 Brockmeier, S.L., Halbur, P.G., Thacker, E.L., 2002. *Porcine Respiratory Disease*  
478 *Complex. Polymicrobial Diseases. American Society of Microbiology.*
- 479 Charrad, M., Ghazzali, N., Boiteau, V., Niknafs, A., 2014. Nbclust: An R package for  
480 determining the relevant number of clusters in a data set. *Journal of Statistical*  
481 *Software* 61, 1-36.
- 482 Chen, Y., Wang, L., Li, F., Du, B., Choo, K.K.R., Hassan, H., Qin, W., 2017. Air quality  
483 data clustering using EPLS method. *Information Fusion* 36, 225-232.

484 Cleveland, R.B., Cleveland, W.S., McRae, J.E., Terpenning, I., 1990. STL: A seasonal-  
485 trend decomposition procedure based on loess. *Journal of Official Statistics* 6, 3-  
486 73.

487 Davenport, F., Funk, C., 2015. Using time series structural characteristics to analyze  
488 grain prices in food insecure countries. *Food Security* 7, 1055-1070.

489 Edwards, D.S., Johnston, A.M., Mead, G.C., 1997. Meat inspection: An overview of  
490 present practices and future trends. *Veterinary Journal* 154, 135-147.

491 Enoe, C., Mousing, J., Schirmer, A.L., Willeberg, P., 2002. Infectious and rearing-system  
492 related risk factors for chronic pleuritis in slaughter pigs. *Preventive Veterinary  
493 Medicine* 54, 337-349.

494 European Community, 2004. Regulation (EC) No 854/2004 of the European Parliament  
495 and of the Council of 29 April 2004 laying down specific rules for the organization  
496 of official controls on products of animal origin intended for human consumption.  
497 OJ L 139,, 206–320.

498 Eze, J.I., Correia-Gomes, C., Borobia-Belsué, J., Tucker, A.W., Sparrow, D., Strachan,  
499 D.W., Gunn, G.J., 2015. Comparison of respiratory disease prevalence among  
500 voluntary monitoring systems for pig health and welfare in the UK. *PLoS ONE* 10.

501 Fablet, C., Dorenlor, V., Eono, F., Eveno, E., Jolly, J.P., Portier, F., Bidan, F., Madec, F.,  
502 Rose, N., 2012. Noninfectious factors associated with pneumonia and pleuritis in  
503 slaughtered pigs from 143 farrow-to-finish pig farms. *Preventive Veterinary  
504 Medicine* 104, 271-280.

505 Fidaner, I.B., Cankorur-Cetinkaya, A., Dikicioglu, D., Kirdar, B., Cemgil, A.T., Oliver,  
506 S.G., 2015. CLUSTERnGO: A user-defined modelling platform for two-stage  
507 clustering of time-series data. *Bioinformatics* 32, 388-397.

508 Hafen, R.P., 2010. *Local regression models: Advancements, applications, and new  
509 methods.* Purdue University.

510 Hennig, C., Meila, M., Murtagh, F., Rocci, R., 2015. *Handbook of Cluster Analysis.* CRC  
511 Press.

512 Leffondré, K., Abrahamowicz, M., Regeasse, A., Hawker, G.A., Badley, E.M., McCusker,  
513 J., Belzile, E., 2004. Statistical measures were proposed for identifying  
514 longitudinal patterns of change in quantitative health indicators. *Journal of Clinical  
515 Epidemiology* 57, 1049-1062.

516 Meriardi, G., Dottori, M., Bonilauri, P., Luppi, A., Gozio, S., Pozzi, P., Spaggiari, B.,  
517 Martelli, P., 2012. Survey of pleuritis and pulmonary lesions in pigs at abattoir  
518 with a focus on the extent of the condition and herd risk factors. *Veterinary  
519 Journal* 193, 234-239.

520 Neumann, E.J., Hall, W.F., Stevenson, M.A., Morris, R.S., Ling Min Than, J., 2014.  
521 Descriptive and temporal analysis of post-mortem lesions recorded in slaughtered  
522 pigs in New Zealand from 2000 to 2010. *New Zealand Veterinary Journal* 62, 110-  
523 116.

524 Niedermeyer, E., Schomer, D.L., da Silva, F.H.L., 2011. Niedermeyer's  
525 Electroencephalography: Basic Principles, Clinical Applications, and Related Fields.  
526 Wolters Kluwer Health/Lippincott Williams & Wilkins.

527 Peng, R., 2008. A Method for Visualizing Multivariate Time Series Data. *Journal of  
528 Statistical Software* 25, (Code Snippet) 1-17.

529 Räsänen, T., Kolehmainen, M., 2009. Feature-based clustering for electricity use time  
530 series data. *Lecture Notes in Computer Science (including subseries Lecture Notes  
531 in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Kuopio, 401-412.

532 Revelle, W., 2015. *psych: Procedures for Psychological, Psychometric, and Personality  
533 Research*. Northwestern University, Evanston, Illinois, USA,.

534 Sanchez-Vazquez, M.J., 2013. *Epidemiological Investigations Utilizing Industry Abattoir  
535 Data – A study in Finishing Pigs*. Utrecht University, Utrecht, The Netherlands.

536 Sanchez-Vazquez, M.J., Nielen, M., Gunn, G.J., Lewis, F.I., 2012. Using seasonal-trend  
537 decomposition based on loess (STL) to explore temporal patterns of pneumonic  
538 lesions in finishing pigs slaughtered in England, 2005-2011. *Preventive Veterinary  
539 Medicine* 104, 65-73.



540 Schuh, M., Köfer, J., Fuchs, K., 2000. Installation of an information feedback system for  
541 control of animal health - Frequency and economical effects of oraan lesions in  
542 slaughter pigs. Wiener Tierarztliche Monatsschrift 87, 40-48.

543 Tabachnick, B.G., Fidell, L.S., 2006. Using Multivariate Statistics (5th Edition).

544 Torgo, R., 2010. Data Mining with R: Learning with Case Studies. Chapman \& Hall/CRC.

545 Vial, F., Reist, M., 2014. Evaluation of Swiss slaughterhouse data for integration in a  
546 syndromic surveillance system. BMC Veterinary Research 10.

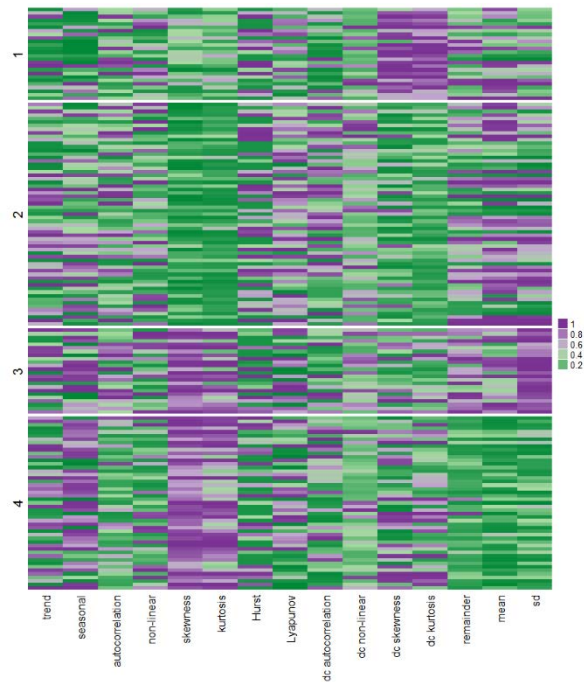
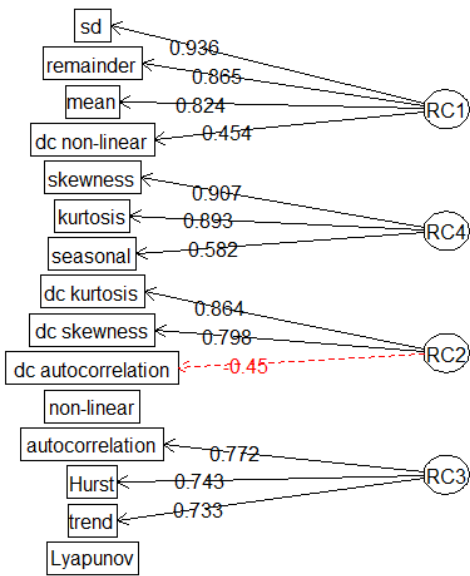
547 Wang, X., Smith-Miles, K., Hyndman, R., 2009. Rule induction for forecasting method  
548 selection: Meta-learning the characteristics of univariate time series.  
549 Neurocomputing 72, 2581-2594.

550 Wang, X., Smith, K., Hyndman, R., 2006. Characteristic-based clustering for time series  
551 data. Data Mining and Knowledge Discovery 13, 335-364.

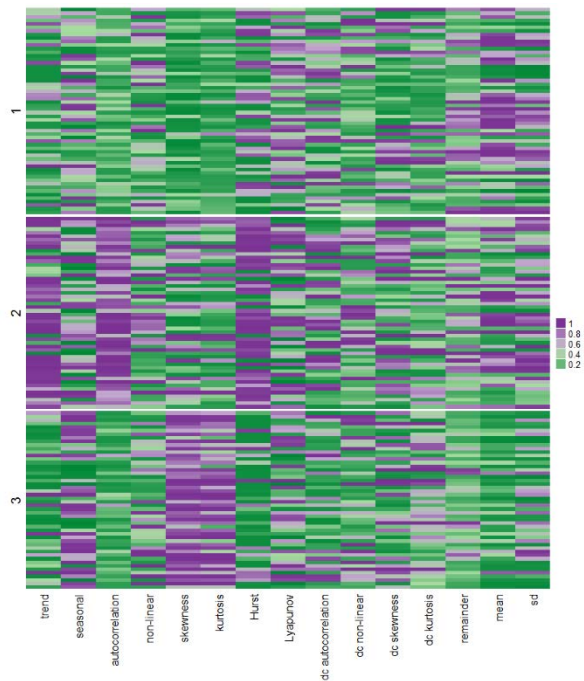
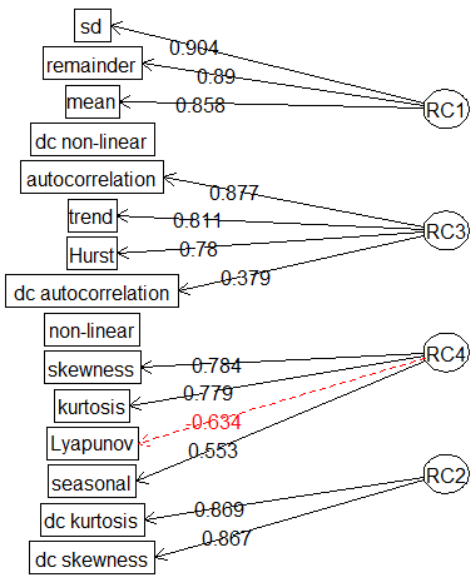
552 Wang, X., Smith, K.A., Hyndman, R.J., 2005. Dimension reduction for clustering time  
553 series using global characteristics. In: Sunderam, V.S., Albada, G.D., Sloot,  
554 P.M.A., Dongarra, J.J. (Eds.), 5th International Conference on Computational  
555 Science - ICCS 2005, Atlanta, GA, 792-795.

556

557

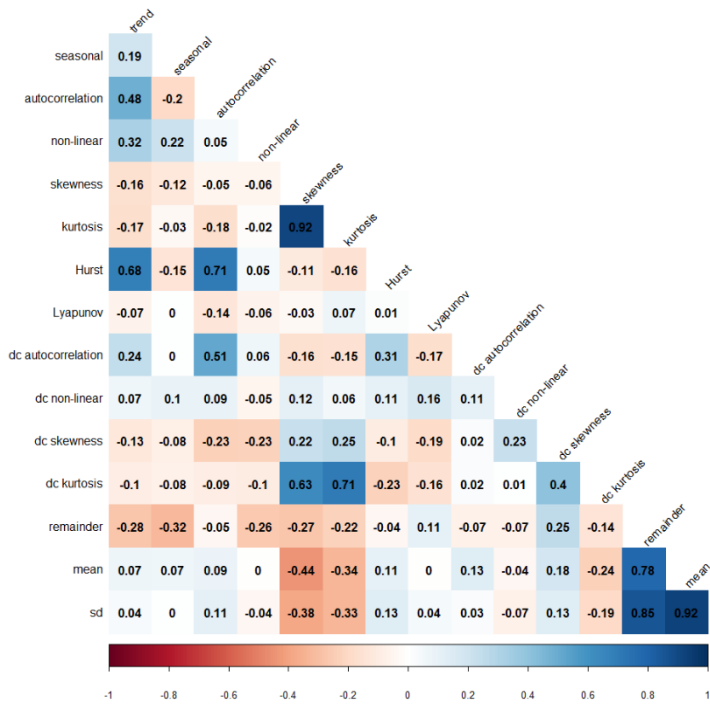


559  
 560 Fig. A1. Factor analysis path diagram of pneumonia (left pane) and allocation to four  
 561 clusters (right pane) in the validation exercise (163 farms). For explanation: see Figure  
 562 3.



563  
 564 Fig. A2. Factor analysis path diagram of pleuritis (left pane) and allocation to three  
 565 clusters (right pane) in the validation exercise (163 farms). For explanation: see Figure  
 566 3.

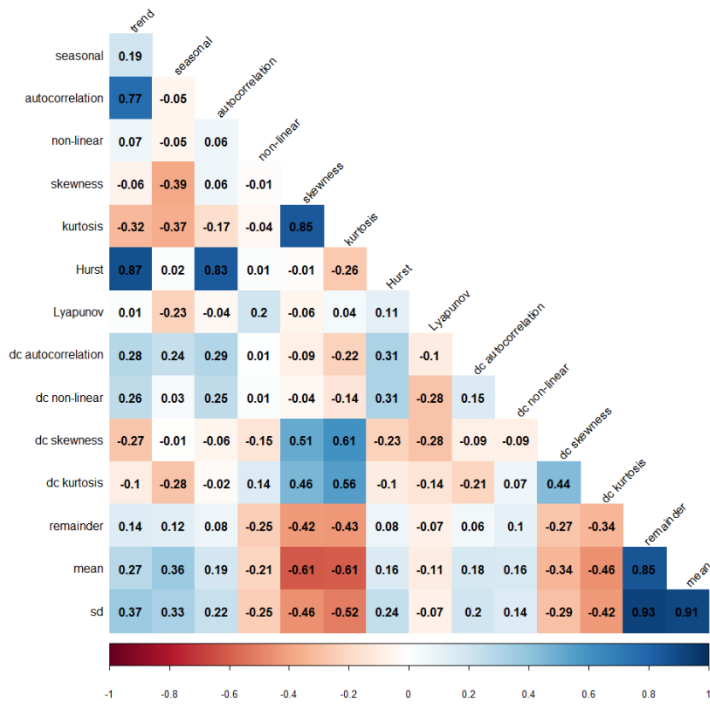
567



568

569 Fig. A3. Correlation matrix for the characteristics of percentage pneumonia

570



571

572 Fig. A4. Correlation matrix for the characteristics of percentage pleuritis

573