

# Clusters of effects in quantile regression models

Gianluca Sottile<sup>a</sup>, Giada Adelfio<sup>a,b,\*</sup>

<sup>a</sup>*Dipartimento di Scienze Economiche Aziendali e Statistiche, University of Palermo, Viale delle Scienze ed. 13, 90128, Palermo, Italy*

<sup>b</sup>*Istituto Nazionale di Geofisica e Vulcanologia, Via Ugo La Malfa, 153 90146, Palermo, Italy*

---

## Abstract

In this paper we propose a new method for finding similarity of effects in a multivariate regression context. Using quantile regression, the effect of each covariate on a response variable is represented as a function of percentiles. Collecting all these curves, describing the effects of each covariate on the response, we could assess if there are covariates with similar effects. Moreover, we provide a flexible algorithm which could be used not only for clustering the coefficient effects of a quantile regression framework, but also for finding clusters of generic curves. We provide also some simulated results and applications on real data, highlighting the flexibility of the proposed approach in several research fields.

*Keywords:* quantile regression coefficients modelling, multivariate analysis, functional data, curves clustering

---

## 1. Introduction

2 In this paper we focus on a new method for classifying effects in general  
3 dependence models. Indeed, a first interest of research could be the comparison  
4 among explanations of different models, that is, if the coefficients associated to  
5 a set of covariates with different responses are different. Another interest could  
6 be to check if there are covariates with similar effects with respect to the same

---

\*Corresponding author

*Email addresses:* [gianluca.sottile@unipa.it](mailto:gianluca.sottile@unipa.it) (Gianluca Sottile),  
[giada.adelfio@unipa.it](mailto:giada.adelfio@unipa.it) (Giada Adelfio)

7 response. Simple t-tests following the ANOVA theory are usually considered to  
8 compare coefficients effects for pooled data, that is, accounting also for some  
9 grouping variable. Extended procedures used to compare regression coefficients  
10 across models (both linear and generalized linear models) are proposed in [5].

11 The novelty of the proposed approach is related to a new perspective of  
12 comparison, focusing not only on single coefficient effect, but on curves effects,  
13 result of a quantile regression fitting.

14 Looking for curve similarity could be a complex issue characterized by sub-  
15 jective choices related to the continuous transformation of observed discrete  
16 data. Here, this problem is handled with the introduction of a new, simple and  
17 efficient procedure, based on a similarity measure between curves. The vari-  
18 ability among curves can be distinguished in two components: phase variability  
19 (removed after the alignment of the curves) and amplitude variability [19].

20 The complex problem of curves clustering is strictly related to the idea of  
21 curves alignment, that is studied in different fields: this is referred to as “curve  
22 registration” in statistics [20, 16], “time warping” in engineering [22] and “struc-  
23 tural averaging” in the context of computing an average curve [12]. A more  
24 general approach is based on the alignment of curves using a target function to  
25 which each one has to be registered with respect to some local features or based  
26 on the minimization of some measure like the average squared distance between  
27 each curve and the target function [20]. [16] used a Procrustes fitting procedure  
28 [9] to provide maximal alignment to the target function, subject to the suitable  
29 smoothness of the transformations. [3] introduced a simple procedure to iden-  
30 tify clusters of multivariate waveforms based on a simultaneous assignation and  
31 alignment procedure. More general methods for curves clustering have been pro-  
32 posed in the literature. [11] introduced a method for finding similarities among  
33 functions by equating the moments between all curves. This problem can be cru-  
34 cial in several contexts. A new approach based on the trimmed K-means Robust  
35 Curve Clustering proposed by [8] is introduced in [2], considering a functional  
36 principal component rotation of data [17]. This approach has been extended  
37 in [4], where the authors focused on finding clusters of multidimensional curves

38 with spatio-temporal structure.

39 All the above mentioned methods have been defined in a slight different  
40 context with respect to the one we consider here. Indeed, the proposed approach  
41 looks for similarities among curves of effects in a quantile regression. These  
42 curves have typically variable trends and different shapes, and the main purpose  
43 is to find effects that are not *significantly* different and could be associated to  
44 covariates belonging to the same cluster, according to a dimensionality reduction  
45 perspective.

46 In general, statistical techniques, aimed at the reduction of huge amounts of  
47 information, are relevant in statistics and synthesis (of objects and variables)  
48 approaches aim to detect the most relevant information for an appropriate in-  
49 terpretation of data.

50 Various methods, combining cluster analysis and the search for a lower-  
51 dimension representation, have been also proposed in the finite dimensional  
52 setting [21]. More recently, the use of clustering is considered as a preliminary  
53 step for exploring data represented by curves, with a further difficulty associated  
54 to the infinite space dimension of data [10].

55 The paper is organized as it follows: in Section 2 we report the usual notation  
56 of Quantile Regression, together with some recent developments referred to a  
57 parametric approach for coefficient functions. In Section 3 we introduce the  
58 new method for curves clustering starting from a quantile regression model,  
59 together with the algorithm details. In Section 4 simulated results are reported  
60 both for curves of effects in quantile regression and in general waveform context.  
61 Example of applications on real data are reported in Section 5. Section 6 is  
62 devoted to conclusive remarks.

## 63 **2. Quantile regression and recent extensions**

64 The non-normality of the distribution and the presence of outliers suggest the  
65 use of Quantile Regression (QR) approach [13, 14] to investigate the influence of  
66 some covariates on the response. Indeed, although the Ordinary Least Squares

67 (OLS) regression allows to model the average as a measure of synthesis, it does  
68 not take into account the whole shape of distribution of the outcome variable.  
69 This issue is overcome by the QR approach: it aims at estimating the fixed  
70 quantiles of the response variable, using different measures of central tendency  
71 (and statistical dispersion), in order to obtain a more comprehensive analysis  
72 of the relationship between variables. In the specific context, the QR analysis  
73 allows to interpret results also for the tails of the distribution, instead of focusing  
74 just on the “average response”. QR deals with the estimation of conditional  
75 quantile functions for models in which quantiles of the conditional distribution  
76 of the response variable are expressed as functions of observed covariates, and  
77 with respect to the usual OLS, QR also provides more robust estimates. Unlike  
78 the ordinary linear regression, the QR parameter measures the change in a  
79 specified quantile of the response variable produced by one unit change in the  
80 predictor variable. This allows to compare how some percentiles of the variable  
81 of interest may be more affected by certain subject characteristics than other  
82 percentiles.

83 In [6], the authors suggest to adopt a parametric model for the coefficient  
84 function of a quantile regression. They refer to this estimation approach as  
85 quantile regression coefficients modelling (QRCM). The QRCM method has  
86 been also implemented in the R package `qrcm` [15, 7].

Conversely to standard quantile regression which works in a quantile-by-  
quantile fashion, in the QRCM framework different quantiles are estimated one  
at the time. This modelling approach facilitates estimation, inference, and  
interpretation of the results, and generally yields a gain in terms of efficiency.  
More in the detail, given a response variable  $y$  and a set of  $q$ -covariates  $\mathbf{x}$ , the  
coefficients  $\boldsymbol{\beta}(p)$  are defined as functions of  $p \in (0, 1)$  (that is the vector of  
percentiles), depending on a finite-dimensional parameter  $\boldsymbol{\theta}$ ,

$$\boldsymbol{\beta}(p \mid \boldsymbol{\theta}) = \boldsymbol{\theta} \mathbf{b}(p),$$

87 where  $\mathbf{b}(p) = [b_1(p), \dots, b_r(p)]^T$  is a set of  $r$  known functions of  $p$  [6]. With this  
88 approach,  $\boldsymbol{\beta}(p)$  is treated as an infinite-dimensional parameter, while the esti-

89 mated coefficients in a standard quantile regression are generally non-smooth  
 90 functions of  $p$  and may suffer from a high volatility that hinders their inter-  
 91 pretability.

92 In a multivariate framework, let  $\mathbf{y} = [y_1, \dots, y_j, \dots, y_m]$  to be a set of  $m$   
 93 response variables, correlated or not, and  $\mathbf{x}$  to be a set of  $q$  covariates. Applying  
 94 the QRCM on each response variable, we estimate the coefficients functions  
 95  $\beta_{1j}(p, \boldsymbol{\theta}), \dots, \beta_{qj}(p, \boldsymbol{\theta})$  over the percentiles. In this paper, starting from the  
 96 QRCM estimation of curve effects, we propose a new algorithm to identify those  
 97 covariates with the same effect on a single response, or, similarly, to identify the  
 98 responses that are related by similar effect of a given covariate. In a generic  
 99 framework, we investigate the similarities among  $n$  general curves, parametrized  
 100 by  $\beta_i(p)$ ,  $i = 1, \dots, n$ .

### 101 3. The proposed clustering method

102 The clustering approach proposed in this paper is based on a new dissimi-  
 103 larity measures based both on shape and distance. More in the detail, we define  
 104 a new dissimilarity measure, based on two measures accounting both for the  
 105 shape and for the distance.

106 Let  $\beta_i(p)$  be the coefficient function approximated by a spline function  $s_i(p)$ ,  
 107 for  $p = 1, \dots, N_p$ ,  $i = 1, \dots, n$ . Considering two different curves  $\beta_i(p)$  and  $\beta_{i'}(p)$   
 108 with  $i \neq i'$ , we define

$$d_{\text{shape}}^{ii'}(p) = I(\text{sign}(s_i''(p)) \times \text{sign}(s_{i'}''(p)) = 1)$$

$$d_{\text{distance}}^{ii'}(p) = I(|\beta_i(p) - \beta_{i'}(p)| \leq f(\alpha, \text{dist}(p))),$$

109 where  $s_i''(\cdot)$  is the second derivative of  $\beta_i(\cdot)$  and  $f(\cdot, \cdot)$  is a cut-off function,  
 110 that depends on  $\alpha$ , a probability value, and  $\text{dist}(p)$ , that is the vector of the  
 111 distances between all the pairs of curves for each value of  $p$ . Therefore, computed  
 112 the distribution of  $\text{dist}(p)$  for each value of  $p$ , the cut-off function selects the  
 113 corresponding  $\alpha$ -th percentile vector.

Therefore, the proposed dissimilarity measure between two curves is defined

as:

$$d^{ii'} = 1 - \frac{1}{N_p} \sum_{p=1}^{N_p} \left[ d_{\text{shape}}^{ii'}(p) \cdot d_{\text{distance}}^{ii'}(p) \right] \quad (1)$$

114 In the proposed approach, the new dissimilarity measure is used to define  
115 a dissimilarity matrix, useful for the application of a hierarchical clustering  
116 method. The proposed procedure has been implemented in the forthcoming  
117 R package `clustEff` that develops some very flexible functions, that allow the  
118 user to make some starting choices. For instance, the  $\alpha$ -level has a central role  
119 for finding homogeneous clusters and its choice can depend on the aim of the  
120 analysis. Fixing an  $\alpha$ -level too small or too big could provide inhomogeneous  
121 clusters. The median is strongly suggested in waveform clustering, while the  
122 first quantile is preferable in clustering of effects. This, of course, could influence  
123 results, but at the same while has the advantage of making the user free to fix  
124 starting conditions according to his/her analysis purpose.

### 125 *3.1. Choice of the number of clusters*

126 In any clustering algorithm, one of the key aspects is the choice of the number  
127 of clusters. In our approach, we deal with this point according to the reference  
128 framework, to provide a classification tool that is both very flexible and could  
129 be used also in different contexts.

130 In a quantile regression context, where the purpose could be to find clusters  
131 of curve effects, we choose the optimal number of clusters (say  $k^*$ ) basing on the  
132 confidence bands of curve. In particular, starting from each partition of curves  
133 in  $k$  clusters and their estimated confidence bands, we build the average band.  
134 Then, we compute the proportion of curves that are outside the average band  
135 (say  $\pi_{out}^k$ ,  $k = 1, \dots, K \leq n$ ). The value of  $k^*$  is identified by that partition for  
136 which  $\pi_{out}^k$  is minimized.

137 Anyway, the proposed approach, based on the dissimilarity measure defined  
138 in (1), could be also an useful tool for clustering of time-dependent signals,  
139 usually analysed in functional data analysis (FDA). The nature of these curves  
140 are different from the one of the effects in a QR. Indeed, in FDA clustering,

141 signals are often zero mean, and with high time-dependent variance. Therefore,  
142 the criterion for the choice of the optimal  $k^*$  can not be the same. In particular,  
143 in waveform clustering framework, we look for the relative distances (say  $\text{dist}_{\text{rel}}^k$ ,  
144  $k = 1, \dots, K$ ) between curves belonging to the same cluster and their centroid.  
145 Then,  $k^*$  is identified by that partition for which the average distance  $\text{dist}_{\text{rel}}^k$  is  
146 minimized.

### 147 3.2. Steps of the Algorithm

148 The main steps of the algorithm are summarized as following:

149 **Step 1.** fixed the  $\alpha$ -level and calculated all the possible distances between the pairs  
150 of curves for each percentile (i.e.  $\text{dist}(p)$ ), the cut-off function selects the  
151 percentile of the distribution of  $\text{dist}(p)$  used in  $d_{\text{distance}}^{i'}(p)$ ;

152 **Step 2.** according to the measure in (1), the dissimilarity matrix is calculated;

153 **Step 3.** applying a hierarchical clustering algorithm a dendrogram is obtained;

154 **Step 4.** if the number of clusters is not fixed, the optimal number is obtained as  
155 in Section 3.1;

156 **Step 5.** after selecting the number of clusters, the mean curves are calculated  
157 within each cluster.

158 The `clustEff` package provides not only the main function that performs the  
159 proposed algorithm, but also a summary and different graphical tools.

160 On the basis of several applications and simulated results, partially here  
161 reported, we can conclude that the algorithm seems to be very stable and fast  
162 in the computation.

## 163 4. Simulation study

164 In this section, we report simulated results for proving the validity of the  
165 proposed approach for cluster of curves, both referring to curves of effects in a  
166 quantile regression and to general waveforms.

167 *4.1. Clusters of effects*

Let us consider a multivariate scenario in which the quantile function is simulated as

$$Q(p | \mathbf{x}, \boldsymbol{\theta}) = \beta_0(p | \boldsymbol{\theta}) + \beta_1(p | \boldsymbol{\theta})x_1 + \cdots + \beta_q(p | \boldsymbol{\theta})x_q,$$

168 where  $x_1, x_2, \dots, x_q$  are independent  $\mathbb{U}(0, 5)$  variables and  $p \in \mathbb{U}(0, 1)$ . In the  
 169 first simulation scenario, the intercept is modelled as a quantile normal distri-  
 170 bution function ( $\phi$ ) for its flexibility. Other choices, as suggested in the original  
 171 paper of [6], could be also considered. We use  $q = 2$  covariates and define three  
 172 groups of quantile functions

$$Q_1(p | \mathbf{x}, \boldsymbol{\theta}) = (1 + \phi(p)) + (.5 + .5p + p^2 + 2p^3)x_1 + (.5 + 2p + p^2 + .5p^3)x_2,$$

$$Q_2(p | \mathbf{x}, \boldsymbol{\theta}) = (1 + \phi(p)) + (-3 + .5p + p^2 + .5p^3)x_1 + (-1.5 - p - .5p^2 + p^3)x_2,$$

$$Q_3(p | \mathbf{x}, \boldsymbol{\theta}) = (1 + \phi(p)) + (.3 - .5p - p^2 + 2p^3)x_1 + (-.5 + p - .5p^2 - p^3)x_2,$$

173 Ten response variables are generated for each quantile function ( $Q_1, Q_2, Q_3$ ).  
 174 Applying the QRCM method to these response variables, we obtained the 30  
 175 coefficients curves, namely curves effect, and their lower and upper bounds,  
 176 useful to select the optimal number of clusters, for both covariates.

177 The `clustEff` algorithm is able to select the correct number of clusters and  
 178 to discriminate the 30 curves effect. In Fig. 1, the curves for both covariates  
 179 are represented in the three clusters, and in Table 1 results are summarized.  
 180 In Table 1 average cluster distances and silhouette widths within clusters are  
 181 reported. The first measure highlights the closeness of curves with respect to the  
 182 mean curve of each cluster, in particular the smaller is this value the closer are  
 183 the curves. The silhouette value is a measure to assess the cohesion of each curve  
 184 to its own cluster compared to other clusters [18]. In particular, observations  
 185 with a large silhouette (almost 1) are very well clustered; a small silhouette  
 186 (around 0) means that there would be some observation that lies between two  
 187 clusters, and negative silhouette means that there are observations probably  
 188 placed in the wrong cluster.

189 Starting from the simulation here reported, results show a valid clustering  
 190 of curves, since silhouette widths are all greater than 0, and in particular, for



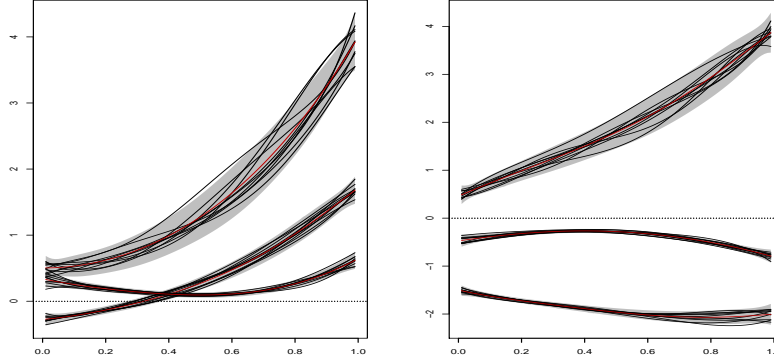


Figure 1: Left and Right panels show the 30 curves clustered in 3 clusters for the first and the second covariate, respectively, after applying the proposed algorithm. Red solid line is the mean curve and dashed red lines are the mean lower and upper bands within each cluster.

Table 1: Results of clustering in correspondence of the two covariates, summarized in terms of average cluster distance (ACD) and silhouette width (SW) within clusters.

	$x_1$		$x_2$	
	ACD	SW	ACD	SW
Cluster 1	.41	.39	.50	.27
Cluster 2	.34	.63	.33	.59
Cluster 3	.45	.77	.27	.86

191 clusters 2 and 3 these values are greater than 0.5. Moreover, all the average  
 192 cluster distances are lower than or equal to 0.5.

#### 193 4.2. Curves clustering

194 Fig. 2 shows 30 curves where 10 of them are obtained from the function  
 195  $f(x) = \sin(3\pi x)$ , 13 from  $g(x) = \cos(3\pi x)$  and 5 from  $h(x) = \sin(3\pi x) \cos(\pi x)$   
 196 evaluated in a grid of size 1000; a  $\mathbb{N}(0, \sigma_t^2)$ -distributed error is added, with  $\sigma_t^2$  a  
 197 variance function defined by segmented relations with multiple change-points.  
 198 Two outlying curves from  $l(x) = 0$  are added, such that they are not pointwise  
 199 outlier at any coordinate. The proposed clustering methods is applied to the 30

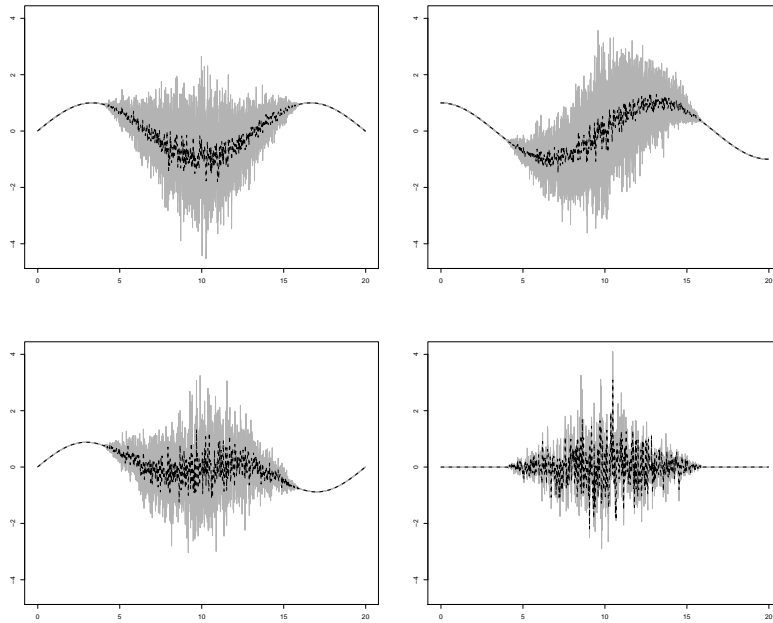


Figure 2: The 30 curves divided in the 4 groups.

200 curves. The applied procedure finds the three clusters  $f(x)$ ,  $g(x)$  and  $h(x)$  and  
 201 also identifies the two outlying curves as a fourth cluster, as reported in Fig. 3.

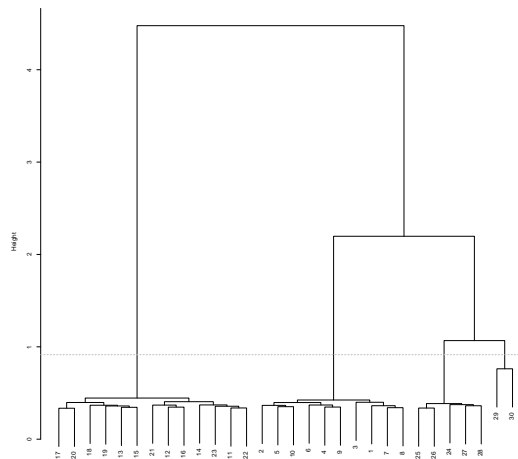


Figure 3: Dendrogram of the clustering algorithm applied in a functional data framework.

202

203 The average distances within the first two clusters is approximately 0.04  
 204 and the individual silhouette width is around 0.41. These results confirm the  
 205 good performance of the proposed method in terms of homogeneity of the found  
 206 clusters and proximity between curves.

## 207 5. Examples of application of the clusteEff algorithm on real data

208 In this section, we apply the proposed clustering algorithm to three different  
 209 real data, in order to show the flexibility of the proposed method and its wide  
 210 spectrum of application.

### 211 5.1. Dataset 1

212 The first analysed dataset consists of 2372 earthquakes located in Italy by  
 213 the INGV (Istituto Nazionale di Geofisica e Vulcanologia) seismic network from

214 2012 to 2016, with local magnitude greater than 2.5. The selected time interval,  
215 as well as the minimum magnitude, have been chosen in order to have a catalogue  
216 as homogeneous as possible. Each seismic event is uniquely identified with a  
217 sequential numeric (ID). For each event Latitude (lat), Longitude (lon) and  
218 Hypocentral Depth (depth), uniquely define the hypocenter position in space.

219 The precision and accuracy of their estimates is strongly influenced by the  
220 quality of the data and the geometry of the stations that recorded the event. In  
221 this application, the following variables are further considered:

- 222 • Magnitude (mag): measure of the magnitude of the earthquake;
- 223 • Magnitude uncertainty (errM): uncertainty about the magnitude of the  
224 earthquake;
- 225 • Hypocentral uncertainty (errZ): uncertainty about the depth hypocenter;
- 226 • Epicentral uncertainty (errH): uncertainty about the depth epicentre;
- 227 • Gap azimuth (gap): a synthetic parameter of the geometry of the stations  
228 in relation to the epicentre; it expresses the maximum angle between two  
229 consecutive stations placing the epicentre to the vertex of the angle. High  
230 values of the azimuthal gap, severely affect the quality of the hypocenter  
231 location. For values higher than  $180^\circ$ , i.e. external seismic event from the  
232 monitoring network, the localization errors can be very high or the event  
233 can not be allocable;
- 234 • Distance from the nearest station (mDst): is the minimum distance be-  
235 tween the epicentre and stations. In particular for shallow earthquakes,  
236 this distance should be sufficiently small. If there is not at least one station  
237 close enough to the epicentre, the determination of depth hypocenter can  
238 be extremely difficult or even impossible. In Figure 4 minimum distances  
239 between epicentres and stations are shown;
- 240 • Root Mean Square (rms): the standard deviation between the arrival times  
241 of seismic waves estimated automatically or manually (experimental) and

242 theoretical ones determined on the basis of a velocity model of wave prop-  
243 agation. This variable is therefore a measure of the quality of the location;

- 244 • Number of stations that recorded the event (nSt): it is the number of  
245 stations used in the localization process. This number is heavily influenced  
246 by the magnitude of the event and strongly influences the accuracy of the  
247 location.

248 Starting from all these variables, we could identify a set of dependent variables  
249 (mag, errM, depth, lon, lat, errZ, errH) and a set of independent variables (gap,  
250 mDst, rms, nSt).

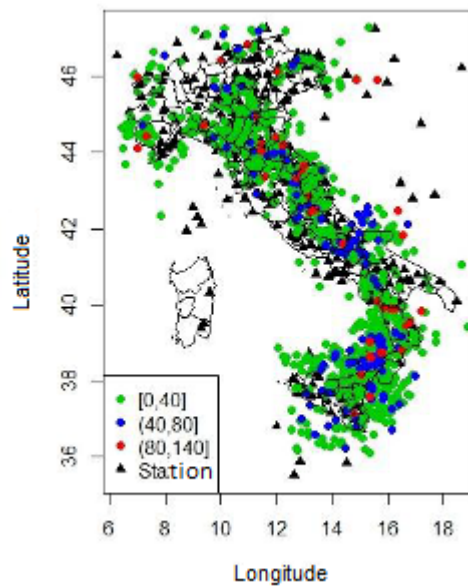


Figure 4: Min. distances (in km) between epicentres and stations (black triangles)

251 One of the main purposes of this analysis is to find some kind of depen-  
252 dence among these sets of variables and, particularly, to identify some clusters  
253 of response variables, conditioning to covariates. Indeed, we look for clusters  
254 of dependent variables after estimating different multiple quantile regressions.  
255 Clustering of effects on different responses, for a fixed covariate, could identify

256 latent relationships among dependent variables. In table 2, we report the corre-  
 257 lation matrix between pairs of variables. As expected, some well known positive  
 258 correlations (errH-gap, errZ-mDst, errH-rms, errZ-rms, errM-rms) and negative  
 259 correlations (errH-nSt, errZ-nSt, gap-nSt) are shown.

Table 2: Correlation matrix between dependent and independent variables. The independent variables are: mag=magnitude, errM=magnitude error, depth, lon=longitude, lat=latitude, errZ=hypocentral uncertainty, errH=Epicentral uncertainty. The dependent variables are: gap=Gap azimuth, dst=distance of the epicentre from the nearest station, rms, nSt=number of stations that recorded the earthquake.

	mag	errM	depth	lon	lat	errZ	errH	gap	dst	rms	nSt
errM	0.03										
depth	0.13	-0.00									
lon	0.01	-0.08	0.33								
lat	-0.01	0.09	-0.38	-0.79							
errZ	0.03	-0.05	0.38	0.19	-0.36						
errH	0.04	-0.12	0.63	0.28	-0.41	0.52					
gap	-0.00	-0.10	0.18	0.20	-0.33	0.32	0.59				
dst	0.15	-0.10	0.31	0.21	-0.38	0.40	0.56	0.61			
rms	0.03	0.01	-0.01	0.03	-0.09	0.15	0.28	0.08	0.11		
nSt	0.52	0.19	0.02	-0.14	0.24	-0.13	-0.21	-0.30	-0.09	0.06	

260 We model the intercept using  $\phi(p)$ , the quantile normal distribution function,  
 261 while the coefficients associated to the covariates are described by a shifted  
 262 Legendre polynomial up to the third degree (e.g., 1), that is, an orthogonal  
 263 polynomial in  $(0, 1)$  used to define a flexible model. In Fig. 5 we report the  
 264 clusters of the dependent variables conditioned to the variables Gap azimuth  
 265 (on the top), RMS (in the middle) and Number of Stations (on the bottom).

266 Conditioning on the Gap azimuth, three clusters of responses are selected:

- 267 1. Magnitude, Magnitude error, Latitude and Hypocentral uncertainty;
- 268 2. Depth and Longitude;
- 269 3. Epicentral uncertainty.

270 In the first cluster, the covariate Gap azimuth has a positive effect on the

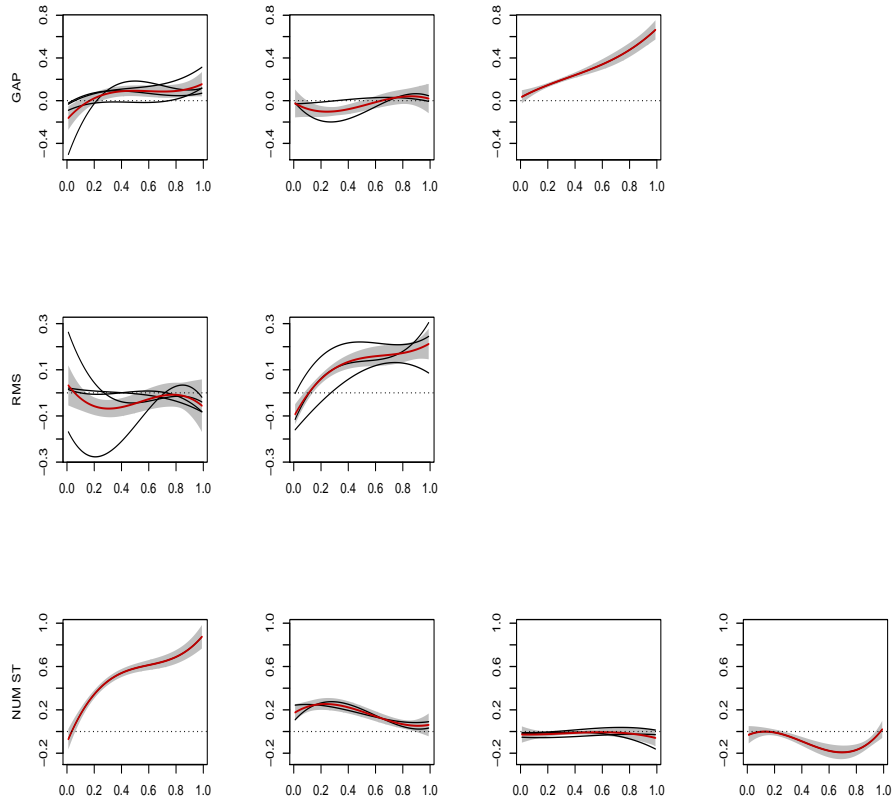


Figure 5: Clusters of responses conditioned to the three independent variables (gap, rms and numst) for dataset 1. Red solid line is the mean curve; the shaded areas are identified by the mean lower and upper bands within each cluster.

271 set of responses just for the percentiles greater than .25, that is, the higher the  
272 Gap the higher both the Magnitude and the Magnitude error. This could be  
273 ascribed to different reasons, i.e. by the distribution of the occurred earthquakes  
274 with respect to the stations. Indeed, the maximum magnitude registered in  
275 the considered catalogue has been the earthquake occurred in June, 2013 with  
276 magnitude 5.9, in the North of Italy (Emilia). For this event, the gap was  
277 very high ( $213^\circ$  vs the average almost  $150^\circ$ ), as well its minimum distance (45  
278 km). The Latitude is also related to the geographical distribution of events.  
279 Indeed, the big sequence of earthquakes occurred after the Emilia event has a  
280 big influence on the estimates. For these events the gaps are high, because in the  
281 North of Italy, and in particular, in the Emilia region, historically considered  
282 as a low seismic area, the network station is less dense. On the other hand the  
283 estimated effect of Gap on the magnitude error and Hypocentral uncertainty,  
284 confirms our previous knowledge.

285 In the second cluster, we could observe a negative effect on the responses for  
286 percentiles lower than .5, and positive otherwise. This reflects the distribution  
287 of station again, and the occurrence features of events. For instance, the area  
288 of the Ionian slab (greater Longitude), where events are deep, the network is  
289 denser (and then the Gap was lower).

290 In the last cluster, an increasing positive effect of Gap azimuth on Epicentral  
291 uncertainty has been estimated, confirming again our previous knowledge.

292 Conditioning on the rms, the responses are clustered in two sets:

- 293 1. Magnitude, Magnitude error, Depth and Longitude;
- 294 2. Latitude, Hypocentral and Epicentral uncertainty.

295 In the first cluster, the effect of the covariate rms on the four responses is  
296 constant and negative for percentiles of the distribution between .15 and .51.

297 In the second cluster, rms has a positive effect on Latitude, Hypocentral and  
298 Epicentral uncertainty, conditioned to the percentiles greater than .13.

299 Finally, conditioning on the covariate nSt, we find four clusters of responses  
300 as follows:



- 301 1. Magnitude;
- 302 2. Magnitude error and Longitude;
- 303 3. Depth, Hypocentral and Epicentral uncertainty;
- 304 4. Latitude.

305 In the first cluster, the number of stations has an increasing positive effect  
306 on Magnitude for almost all the distribution.

307 In the second cluster, the number of stations has a decreasing positive effect  
308 on Magnitude error and Longitude.

309 In the third cluster, the covariate number of stations has a low negative effect  
310 on Depth, Hypocentral and Epicentral uncertainty for percentiles between .2 and  
311 .36.

312 In the last cluster, the number of stations has almost always a negative effect  
313 on Latitude, reflecting what we observed with respect to the covariate Gap.

314 In this application, we show an interesting usage of the proposed method,  
315 identifying clusters of dependent variables on the basis of the estimated curves  
316 effect for a better characterization of these dependencies.

## 317 5.2. Dataset 2

318 The data refer to a study carried out in 1988-1991 in the North of Italy, in-  
319 cluding 1053 males and 992 females. The study aims at assessing determinants  
320 of the Inspiration Capacity (IC), a measure of lung's function, among the fol-  
321 lowing nine predictors: age, height, body mass index (bmi), sex, and indicators  
322 for current smoking, occupational exposure, cough, wheezing, and asthma.

323 We model the intercept using  $\log(p)$  and  $\log(1-p)$ , that defines the asym-  
324 metric Logistic distribution used for its flexibility, while the coefficients associ-  
325 ated to the covariates are described by a fifth degree shifted Legendre polyno-  
326 mial. The estimated model is summarized in Table 3, and curves of effects are  
327 represented in Figure 6.

328 We apply the `clustEff` algorithm to the curves of significant coefficients  
329 associated to the predictors in order to look for similar effects of covariates with

Table 3: Summary of the model selected by AIC, for the Inspiratory Capacity data with respect to the nine listed covariates

	1	$\log(p)$	$\log(1-p)$	$slp(p, 5)[1]$	$slp(p, 5)[2]$	$slp(p, 5)[3]$	$slp(p, 5)[4]$	$slp(p, 5)[5]$	P-value
(Intercept)	-0.050	0.353	-0.403	-	-	-	-	-	.000*
age	-0.285	-	-	0.049	0.002	0.012	-0.013	-0.000	.000*
hgt	0.411	-	-	0.085	0.086	0.031	0.024	-0.006	.000*
bmi	0.146	-	-	0.069	-0.052	0.016	-0.019	0.015	.000*
sex	-0.016	-	-	-0.101	0.077	-0.022	0.055	-0.010	.000*
smoker	-0.060	-	-	0.041	-0.012	0.000	-0.023	-0.002	.208
occ_exp	0.085	-	-	0.008	0.020	-0.018	0.021	-0.016	.594
cough	0.018	-	-	-0.005	0.007	-0.028	-0.009	-0.009	.306
wheeze	-0.013	-	-	0.054	0.017	-0.007	0.037	0.010	.048*
asthma	0.007	-	-	0.004	-0.018	0.020	0.010	-0.035	.380
P-value	.000*	.000*	.000*	.000*					

QRCM estimates of  $\theta$ . The corresponding graphical representation is proposed in Figure 6. The last row reports the p-values for the null hypothesis that the coefficient of  $\theta$  is  $\mathbf{0}$  and represents the significance of the components of  $\mathbf{b}(p)$ . The last column is defined analogously and can be interpreted as the significance of a test for a null effect of covariates. The asterisk (\*) denotes significance less than 0.05.

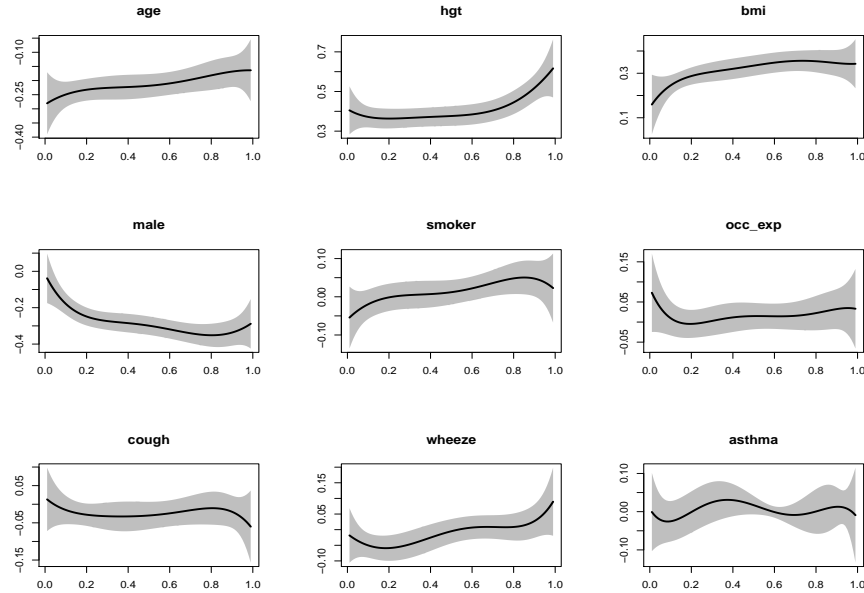


Figure 6: QRCM estimates of  $b(p)$  (see Table 3). Confidence bands are displayed as shades.

330 respect to the Inspiration Capacity response. The application of the proposed  
 331 algorithm provides three clusters.

332 The first cluster consist of age and sex with an average negative effect, sig-  
 333 nificant for all the percentiles.

334 The second cluster is identified by bmi and height, with a positive effect,  
 335 significant for all the percentiles.

336 Finally, in the third cluster we find the variables smoker, occupational ex-  
 337 posure, cough, wheezing, and asthma with not significant average effect.

338 These results are summarized in Figure 7.

339 In this application, we focus on a new perspective of reduction of dimen-  
 340 sionality, applied in a quantile regression context. Indeed, we propose the use  
 341 of the `clustEff` method for finding the main determinants of a quantitative  
 342 response, assuming that we are interested in looking for dependence structures.  
 343 Of course, these results could be more relevant in presence of several regressors,

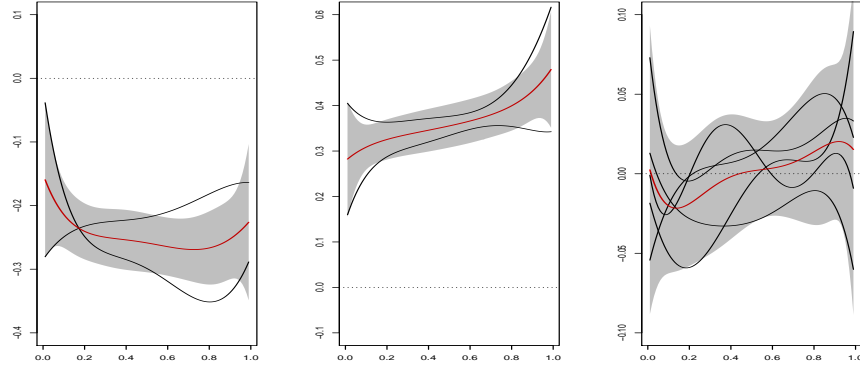


Figure 7: The three clusters obtained applying the `clustEff` algorithm on the estimated quantile regression coefficients of dataset 2. Red line is the mean curve; the shaded areas are identified by the mean lower and upper bands within each cluster. Black lines are the covariates; in the first cluster there are the variables age and sex, in the second cluster bmi and height and in the third cluster smoker, occupational exposure, cough, wheezing and asthma.

344 but we showed this example just for its simplicity of interpretation. Indeed,  
 345 we could observe that covariates are classified in three main groups: the first  
 346 relative to the subject characteristics, the second relative to body features and  
 347 the third that associates the clinical aspects. Therefore, in describing the effect  
 348 of covariates to the response, we interpret the average effect of each cluster, as  
 349 a proxy of a latent characteristic effect that is associated to the covariates of  
 350 that cluster. As drawback, this procedure could have some limitations in terms  
 351 of loss of interpretation, as usual in dimensionality reduction problems.

### 352 5.3. Dataset 3

353 Also this application is reported to show the flexibility of the proposed algo-  
 354 rithm. Indeed, we used the `clustEff` method for waveform clustering, that may  
 355 be considered as an issue of clustering of functional data.

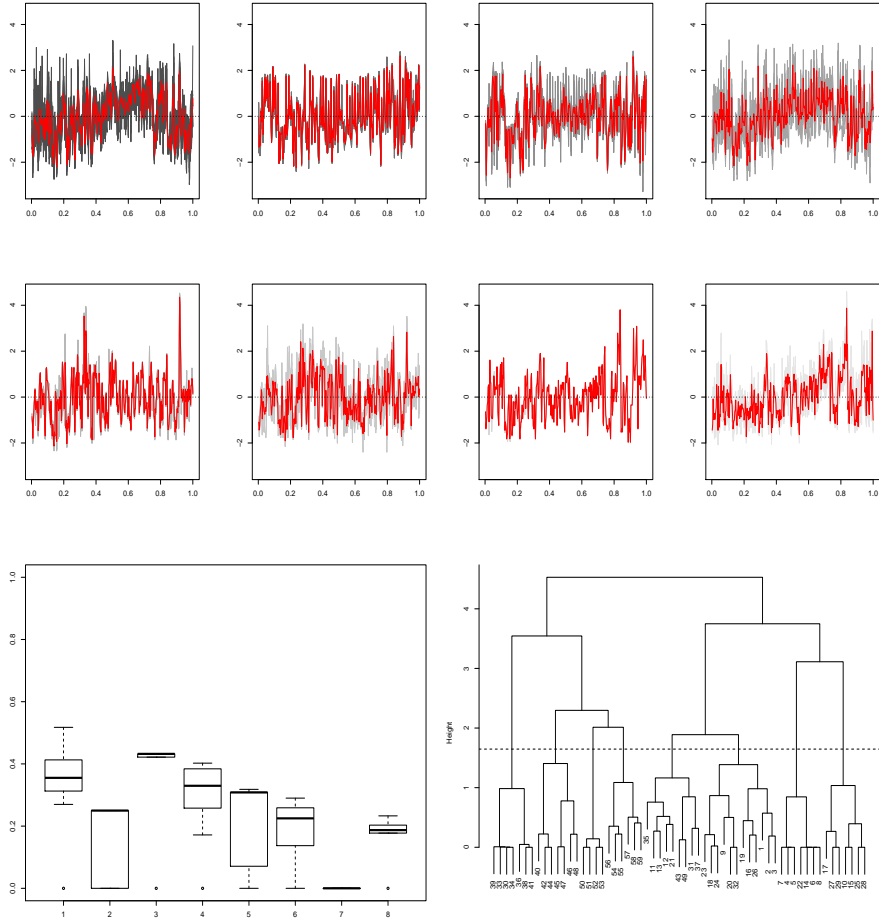


Figure 8: The identified 8 clusters of dataset 3 (on the top): red lines are the mean curves. Boxplot of the average mean distance within each cluster (on the left-bottom). Dendrogram of the clustering algorithm and height level used to cut the tree (on the right-bottom).

## 356 **6. Conclusion**

357 The proposed approach is not just a method for clustering of curves, that  
358 is an important problem in many areas of science, but it can be seen as a  
359 new tool for reduction of dimensionality in dependence model, in particular  
360 in a quantile regression context. Indeed, the proposed approach, based on a  
361 new dissimilarity measure, that accounts both for shape of curves and distance  
362 among them, allows to find similarities among curves that represent the effect  
363 of covariates on (also multivariate) response. The clustering of these curves,  
364 extends the idea of looking for similar effects and, therefore, of covariates in  
365 general dependence models, aimed to a selection perspective.

366 This approach, developed also in a forthcoming R Package, is a very flex-  
367 ible method, that is also very fast in te of computation and user-friendly for  
368 general applications. We applied the proposed algorithm to three different real  
369 data, included an application for generic waveforms in order to provide a wider  
370 spectrum of applications for curves clustering.

## 371 **Acknowledgements**

372 This paper has been partially supported by the national grant of the Italian  
373 Ministry of Education University and Research (MIUR) for the PRIN-2015 pro-  
374 gram (Progetti di ricerca di Rilevante Interesse Nazionale), “Prot. 20157PRZC4  
375 - Research Project Title Complex space-time modelling and functional analysis  
376 for probabilistic forecast of seismic events. PI: Giada Adelfio”

## 377 **References**

- 378 [1] Abramowitz, M. and Stegun, I.A. (1964). Handbook of mathematical func-  
379 tions with formulas, graphs, and mathematical tables. Courier Corporation,  
380 vol. 55.
- 381 [2] Adelfio, G., Chiodi, M., D’Alessandro, A. and Luzio, D. (2011). FPCA al-  
382 gorithm for waveform clustering. Journal of Communication and Computer,  
383 vol. 8(6):494-502. ISSN 1548-7709

- 384 [3] Adelfio, G., Chiodi, M., D'Alessandro, A. and Luzio, D., D'Anna, G.,  
385 Mangano, G. (2012) Simultaneous seismic wave clustering and registration.  
386 Computers & Geosciences, 8(44), 6069.
- 387 [4] Adelfio, G, Di Salvo, F, Chiodi, M. (2016). Space-time FPCA Algorithm  
388 for clustering of multidimensional curves. Proceeding of the 48th Scientific  
389 Meeting of the Italian Statistical Society, Salerno, June 8th June 10th, 2016.  
390 Editors: Monica Pratesi and Cira Pena. ISBN: 9788861970618
- 391 [5] Clogg, C, Petkova, E, Haritou, A. (1995). Statistical Methods for Comparing  
392 Regression Coefficients Between Models. The American Journal of Sociology,  
393 100, 5: 1261-1293
- 394 [6] Frumento, P. and Bottai, M., (2015). Parametric modeling of quantile re-  
395 gression coefficient functions, Biometrics, 72:74-84.
- 396 [7] Frumento, P. (2017). qrcm: Quantile Regression Coefficients Modeling. R  
397 package version 2.0, <https://CRAN.R-project.org/package=qrcm>.
- 398 [8] Garca-Escudero, L. A. and Gordaliza, A. (2005). A proposal for robust curve  
399 clustering, Journal of classification, 22, 185-201.
- 400 [9] Gower, J., 1975. Generalized Procrustes analysis. Psychometrika 40(1):33-  
401 51.
- 402 [10] Jacques J. and Preda C. (2014). Functional data clustering: a survey. Ad-  
403 vances in Data Analysis and Classification, Springer Verlag, 8:231-255
- 404 [11] James, G.M., (2007). Curve alignment by moments. The Annals of Applied  
405 Statistics 1:480-501.
- 406 [12] Kneip, A. and Gasser, T., (1992). Statistical tools to analyze data repre-  
407 senting a sample of curves. Annals of Statistics 20:1266-1305.
- 408 [13] Koenker, R. and Bassett, G., Jr., (1978). Regression quantiles, Economet-  
409 rica, 46:33-50.

- 410 [14] Koenker, R. (2005). *Quantile Regression*, Cambridge University Pres.
- 411 [15] R Core Team (2016). *R: A language and environment for statistical com-*  
412 *puting*. R Foundation for Statistical Computing, Vienna, Austria. URL  
413 <https://www.R-project.org/>.
- 414 [16] Ramsay, J.O. and Li, X., (1998). Curve registration. *Journal of the Royal*  
415 *Statistical Society, Sec. B*, 60:351-363.
- 416 [17] Ramsay, J. O. and Silverman, B. W. (2005). *Functional Data Analysis*.  
417 Springer, New York.
- 418 [18] Rousseeuw, P. J. (1987). Silhouettes: a Graphical Aid to the Interpretation  
419 and Validation of Cluster Analysis. *Computational and Applied Mathematics*.  
420 20: 5365.
- 421 [19] Sangalli, L.M., Secchi, P., Vantini, S. and Veneziani, A., (2009). A Case  
422 Study in Explorative Functional Data Analysis: Geometrical Features of  
423 the Internal Carotid Artery. *Journal of the American Statistical Association*  
424 104(485):37–48.
- 425 [20] Silverman, B.W., (1995). Incorporating parametric effects into functional  
426 principal components analysis. *Journal of the Royal Statistical Society, Sec.*  
427 *B* 57:673-689.
- 428 [21] Vichi, M., Saporta, G. (2009). Clustering and Disjoint Principal Compo-  
429 nent Analysis. *Computational Statistics and Data Analysis*, 53:3194-3208.
- 430 [22] Wang, K. and Grasser, T., (1997). Alignment of curves by dynamic time  
431 warping. *The Annals of Statistics* 25:1251-1276.