



Analysing radon accumulation in the home by flexible M-quantile mixed effect regression

R. Borgoni¹ · A. Carcagni¹ · N. Salvati² · T. Schmid³

© Springer-Verlag GmbH Germany, part of Springer Nature 2018

Abstract

Radon is a noble gas that occurs in nature as a decay product of uranium. Radon is the principal contributor to natural background radiation and is considered to be one of the major leading causes of lung cancer. The main concern revolves around indoor environments where radon accumulates and reaches high concentrations. In this paper, a semiparametric random-effect M-quantile model is introduced to model radon concentration inside a building, and a way to estimate the model within the framework of robust maximum likelihood is presented. Using data collected in a monitoring survey carried out in the Lombardy Region (Italy) in 2003–2004, we investigate the impact of a number of factors, such as geological typologies of the soil and building characteristics, on indoor concentration. The proposed methodology permits the identification of building typologies prone to a high concentration of the pollutant. It is shown how these effects are largely not constant across the entire distribution of indoor radon concentration, making the suggested approach preferable to ordinary regression techniques since high concentrations are usually of concern. Furthermore, we demonstrate how our model provides a natural way of identifying those areas more prone to high concentration, displaying them by thematic maps. Understanding how buildings' characteristics affect indoor concentration is fundamental both for preventing the gas from accumulating in new buildings and for mitigating those situations where the amount of radon detected inside a building is too high and has to be reduced.

Keywords Environmental radioactivity · Building factors · Radon-prone areas · Hierarchical mixed models · Penalised splines · Lombardy region

1 Introduction

Radon (the ²²²Rn isotope) is a noble, radioactive gas naturally occurring as a decay product of uranium. It is a gas without colour or smell that is detectable only by specialised measurement devices and represents the main contributor to natural background radiation. The becquerel (Bq) is the standard international unit for radon activity i.e. the amounts of radioactive material, and radon activity concentration is measured in Bq/m³.

Evidence since the sixteenth century suggests that exposure to elevated concentrations of radon and radon progeny is a potential cause of the high prevalence of lung cancer mortality among miners (Jacobi 1993). However, it is only since the 1970s that human exposure to radon became of general concern. Radon is present both indoors and outdoors. While outdoors, its concentration in the air is diluted to a low level and therefore does not pose any significant health problems. Indoor concentrations are far higher simply because the gas enters into a smaller space and accumulates there. Hence, this paper focusses on indoor radon concentration (IRC), since it is in the indoor environment that radon becomes a serious health concern.

Extensive epidemiological studies (Lubin and Boice 1997; Kreienbrock et al. 2001; Darby et al. 2005; Krewski et al. 2005; Tiefelsdorf 2007) point out that long-term radon exposure in homes determines a remarkable increase in the risk of lung cancer. Nowadays, the International

✉ R. Borgoni
riccardo.borgoni@unimib.it

¹ Università degli studi di Milano - Bicocca, Milan, Italy

² Università di Pisa, Pisa, Italy

³ Freie Universität Berlin, Berlin, Germany

Agency for Research on Cancer (IARC) and the US Environmental Protection Agency (EPA) have categorised radon as a Group 1 and Group A human carcinogen, respectively. Other papers have investigated the impact of radon on the increase in the risk of other cancer typologies (Smith et al. 2007).

To estimate the exposure of people to radon and to identify building typologies and geographical areas more prone to high IRC, monitoring surveys have been implemented in many countries such as in the UK (Green et al. 2002) USA (USEPA 1992) (Smith and Field 2007), Canada (Shi et al. 2006), and Belgium (Cinelli et al. 2011), to mention but a few. These surveys provide geocoded data that are often fundamental for planning remediation activities. Collecting information at several locations, data coming from these monitoring campaigns allow us to account for the multifactorial dependencies of IRC through statistical models that combine a number of explanatory variables. Raised IRC levels can often be traced back to the radon content in the underlying rocks and soils and are detected in dwellings close to the ground. Hence, the geological and lithologic nature of the soil as well as other soil characteristics, such as porosity and permeability, can influence indoor accumulation. Radon exhalation from building materials is another relevant source of radon since many building materials, such as concrete containing alum shale or volcanic tuffs and pozzolana, may have high radium content. The relationship between IRC and geological indicators of high radon potential as well as other structural building factors and materials has long been documented: see Gunby et al. (1993), Gates and Gundersen (1992), Price et al. (1996), Apte et al. (1999), Levesque et al. (1997), Sundal et al. (2004), Shi et al. (2006), Smith and Field (2007), Hunter et al. (2009), Cinelli et al. (2011) amongst others. All these types of effects will be considered further in the paper.

All the papers mentioned so far aimed to assess the influence of various characteristics of buildings on the average IRC. However, modelling a central tendency measure of the conditional distribution, typically the mean, may provide a rather incomplete or even inappropriate picture if the actual interest is in the tail of the distribution as is the case in the presence of outliers, asymmetry or reference concentration values endorsed by law or international recommendations. The conventional regression approach, modelling the conditional expectation of IRC, does not permit this kind of analysis. Instead, this can be obtained using quantile regression, i.e. by fitting a family of robust regression models, each summarising the behaviour at different levels of the IRC conditional distribution. Quantile regression was introduced in the econometrics literature by Koenker and Bassett (1978) and has been extended towards many directions. For instance, Chaudhuri

(1991) proposed locally polynomial quantile regression whereas Koenker et al. (1994) and Bosch et al. (1995) discussed penalty methods for smoothing quantile regression. Geraci and Bottai (2014) investigated linear quantile regression models for clustered/hierarchical data. Since quantile regression enables us to fully describe the conditional distribution, the method has been used in many applications in recent decades. We refer to Yu et al. (2003) and Koenker (2005) for some examples. However, quantile regression has seldom been applied in the context of radon mapping. Exceptions are represented by Borgoni (2011) who adopted a spatial semiparametric model to define a conditional quantile regression model. Furthermore, Fontanella et al. (2015) and Sarra et al. (2016) proposed a spatial quantile hierarchical Bayesian model in this context.

An alternative to quantile regression is the M-quantile regression introduced by Breckling and Chambers (1988) to integrate expectile (Newey and Powell 1987) and quantile regression within a unique paradigm based on a 'quantile-like' generalisation of regression defined via influence functions (M-regression). Tzavidis et al. (2016) extended the M-quantile regression by including random effects in order to consider the hierarchical structure in the data, whereas Alfó et al. (2017) proposed a finite mixture of M-quantile regressions with discrete random coefficients; the discrete distribution of the latter can be interpreted as a nonparametric estimate of an unspecific continuous distribution. Although M-quantile and quantile regressions cannot be directly compared as they target different location parameters, both approaches try to model location parameters that are related to the same part of the conditional distribution (Jones 1994). Why should the potential data analyst consider M-quantile regression when the main advantage of quantile regression is the more intuitive interpretation? M-quantile regression models are more flexible, in particular they allow for robustness in exchange for efficiency in inference by tuning a suitable constant of the influence function (see Sect. 3). The option to select different continuous influence functions in an M-quantile regression—in contrast to the absolute value function in a quantile regression—can offer additional computational stability.

As far as the statistical methodology is concerned, this paper extends the work by Tzavidis et al. (2016) to a random effect semiparametric M-quantile model which is able to account both for the characteristics of the soil and for the material and architectural structure of a building. It is well known, however, that radon dynamics are spatially structured due to a number of causes that may affect IRC on a local and on a large-scale over and above the available secondary information obtained via administered surveys or measurement campaigns. For this reason, we extend the

158 model to include a flexible component that is able to grasp
159 this spatial effect. In particular, the proposed M-quantile
160 model incorporates the spatial information (locations) of
161 the data through a spline component in the linear predictor
162 of the model, and therefore does not rely on any structural
163 assumptions in the error terms. This is particularly relevant
164 when, as in the case study presented later in the paper,
165 geographically referenced measures have to be spatialised
166 to produce maps.

167 The paper is structured as follows. Section 2 presents the
168 data and the model we propose in Sect. 3. Section 4 shows
169 the results obtained applying the suggested random effect
170 semiparametric M-quantile model to indoor radon data. In
171 particular, we explain how the suggested model permits us
172 to map the phenomenon of interest across space and
173 identify those areas more prone to high concentration and
174 to estimate the impact of potential determinants on the
175 pollutant concentrations. Concluding remarks are presented
176 in Sect. 5.

177 2 Data description

178 The data used in the present paper come from an indoor
179 radon gas monitoring survey implemented by the Agency
180 for Environmental Protection (ARPA) from 2003 to 2005
181 in the Lombardy Region (northern Italy). With a surface
182 area of 23,800 km² Lombardy is the fourth largest and
183 most populated region in Italy, with about 10,000,000
184 inhabitants according to the last census in 2011, which
185 corresponds to about 20% of the entire Italian population.
186 A national survey conducted by the National Health Ser-
187 vice from 1989 to 1994 already indicated that Lombardy is
188 exposed to high values of IRC. This survey pointed out that
189 the IRC in Lombardy was 116 Bq/m³ on average, and is
190 therefore higher than the national average of 70 Bq/m³.
191 Assessing the spatial variability of IRC and the population
192 exposure to this gas is a prominent environmental and
193 health-related issue in this part of the country.

194 In this paper, the problem of high IRC in dwellings is
195 investigated by considering a sample of 900 measures of
196 IRC collected throughout the regional territory for which a
197 complete record of all relevant building characteristics
198 were available. Figure 1 shows the locations of the mea-
199 surement points (indicated by crosses) and the study region
200 expressed in UTM projection.

201 Long-term measurements was obtained using CR-39
202 trace detectors that were positioned in dwellings for
203 12 months. The dosimeters were changed after approxi-
204 mately 6 months and the year-long average of the two
205 semester values is considered in this paper, weighting the
206 two one-semester measurements by the actual time of
207 exposure of each detector. The average concentration is

around 118 Bq/m³ (sd 136 Bq/m³) ranging from a mini- 208
209 mum of 12.5 Bq/m³ to a maximum 1762.5 Bq/m³. As
210 shown in Fig. 2, the IRC distribution is strongly asym-
211 metric with a number of potential outliers, in line with
212 other studies (Nero et al. 1986, amongst others).

213 A questionnaire were also administered to dwellers of
214 each sampled unit to collect other information about the
215 building and the rooms in addition to the IRC. IRC mea-
216 surements and building information were then combined in
217 a single dataset. In the following, we focussed on factors
218 that are expected to affect IRC, such as the wall material
219 (stone versus other materials such as lateritious and hollow
220 brick), the presence of an air conditioning system, the type
221 of connection with the soil (i.e. whether the building is in
222 direct contact with the ground or a basement/crawlspace is
223 present), the type of building (detached vs. non-detached),
224 the year of construction or last renovation (before or after
225 1990) and the floor material (marble or granite versus other
226 materials). In Table 1 some summary statistics of IRC
227 conditioned to these building factors are shown. We
228 observe that different house characteristics impact differ-
229 ently on the IRC level. For instance, the differences
230 between the 20th percentiles and 80th percentiles of IRC
231 measured in dwellings with a marble-granite floor and in
232 dwellings with an other-material floor are 4.9 Bq/m³ (20th
233 percentile) and 14.2 Bq/m³ (80th percentile), respectively.
234 When comparing the 20th percentiles and 80th percentiles
235 of IRC measured in dwellings with stone walls versus other
236 material walls, the differences are 11.1 Bq/m³ and
237 62.1 Bq/m³, respectively, with a much more pronounced
238 spread between higher quantiles. Estimating the quantiles
239 of the IRC conditional distribution, given the covariates, is
240 worth pursuing.

241 The composition of the soil on which a building is
242 located is another important feature that can affect the IRC,
243 since the concentration of uranium and radium varies
244 depending on the rock lithology. Hence, it is expected that
245 higher concentration levels tend to occur in particular
246 geological areas. The geological composition of Lombardy
247 varies extremely with regards to the lithological and soil
248 typologies. In order to derive this information the data were
249 linked to a geo-lithologic map on a scale of 1:250,000
250 (Borgoni et al. 2011) that partitions the territory of the
251 Lombardy into 11 geological classes (see Fig. 3). Since the
252 measurement points were geo-referenced, it was possible to
253 assign one of the 11 types to each of them.

254 Figure 4 shows the boxplot of IRC for each geological
255 class. The dashed lines in the figure connect the 20th and
256 80th quantiles of IRC conditioned to different geo-litho-
257 logic classes. We observe that the quantiles change con-
258 siderably between the geo-lithologic classes.

259 Finally, high IRC can also be found in areas with low
260 radium levels, especially when fractured rocks or intensive

Fig. 1 Sampling locations and the study area

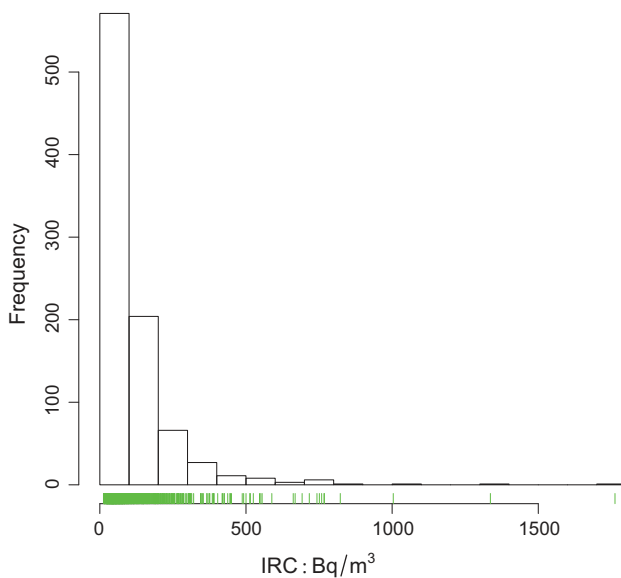
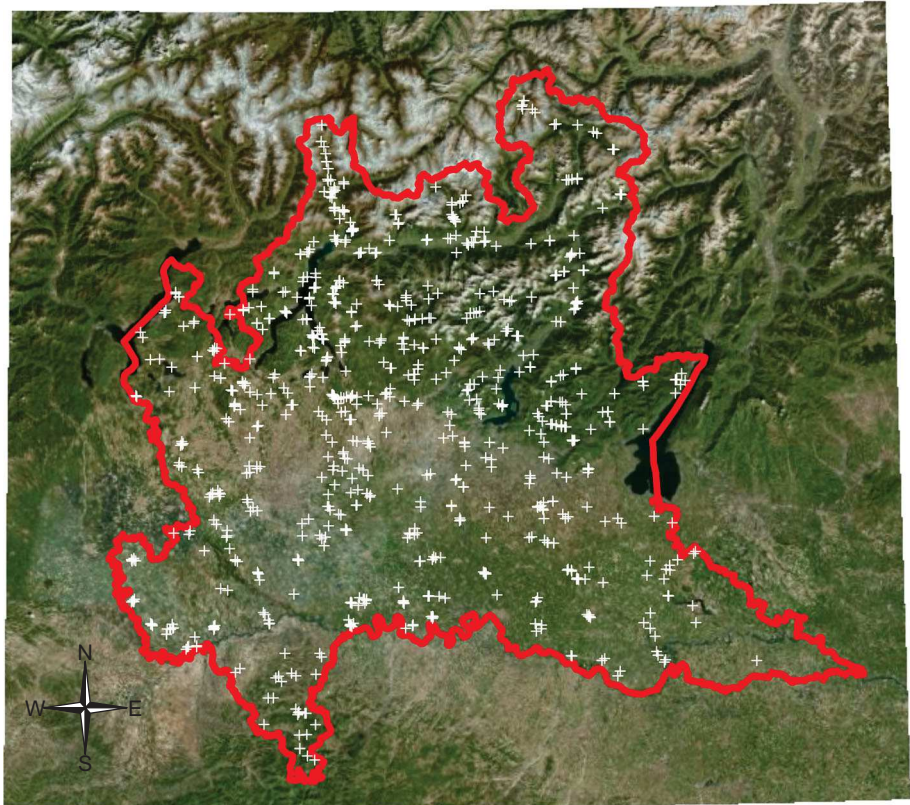


Fig. 2 Sampling distribution of IRC

261 tectonic frameworks are present. This may be due to the
 262 presence of faults, which may foster the gas to seep up
 263 from deeper origins and enter into homes. To assess the
 264 effect of being in the proximity of a fault, we calculated the
 265 distance of each sampling point x to a tectonic fragment
 266 A . Tectonic fragment cartography was available from a
 267 shape file where each fault is geocoded in a vector format

through a set of nodes. Hence, the distance has been calculated as

$$d(x, A) = \min_{s \in A} \|x - s\| \tag{1}$$

as suggested by Foxall and Baddeley (2002).

3 A random effect semiparametric M-quantile model for IRC

The M-quantile of order q , $MQ_y(q|\mathbf{X}; \psi)$, for the conditional density of an outcome variable y given auxiliary variables \mathbf{X} , $f(y|\mathbf{X})$, is defined by Breckling and Chambers (1988) as the solution of the integral equation $\int \psi_q \{y - MQ_y(q|\mathbf{X}; \psi)\} f(y|\mathbf{X}) dy = 0$. Here, ψ_q is the derivative of an asymmetric loss function ρ_q , called the (asymmetric) influence function. In particular, (\mathbf{x}_i^T, y_i) , $i = 1, \dots, n$, denotes n observations of a random sample, y_i is the outcome variable and \mathbf{x}_i^T are the p -vectors of the covariates \mathbf{X} . A linear M-quantile regression model of y_i given the auxiliary variables \mathbf{x}_i is given by

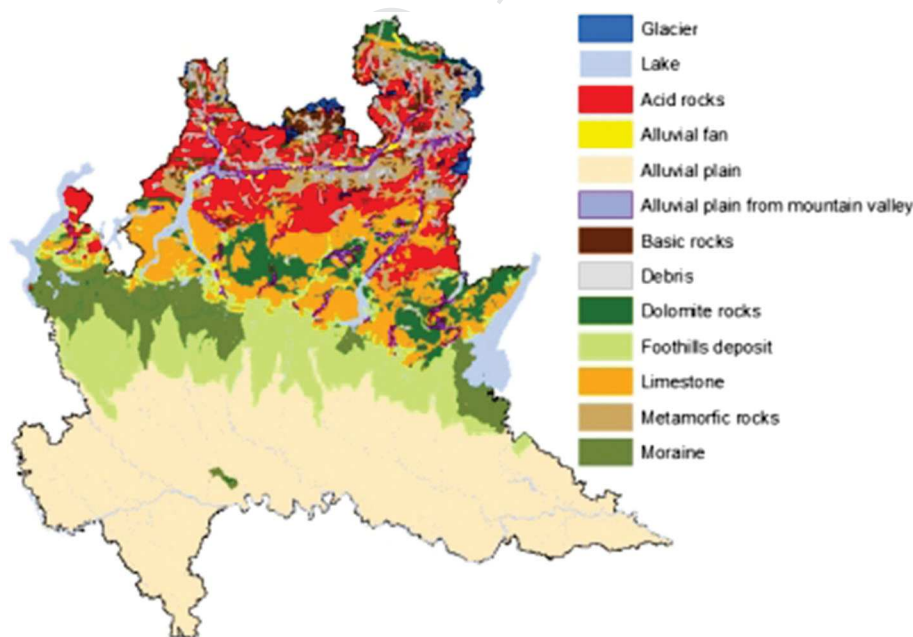
$$MQ_{y_i}(q|\mathbf{x}_i; \psi) = \mathbf{x}_i^T \boldsymbol{\beta}_q,$$

and $\boldsymbol{\beta}_q$ can be estimated minimising

Table 1 IRC summary statistics by dwelling characteristics

	N	Mean	SD	Min	Q20	Median	Q80	Max
<i>Wall material</i>								
Other	788	110.4	114.3	12.5	40.3	73.2	150.7	1003.9
Stone	112	171.1	233.3	16.0	51.2	98.7	212.8	1762.5
<i>Conditioning system</i>								
No	840	120.2	138.8	12.5	41.9	76.7	164.7	1762.5
Yes	60	85.4	84.5	19.9	40.6	65.7	95.7	555.2
<i>Connection with the ground</i>								
In contact	352	135.9	143.0	13.8	45.0	84.6	195.3	1003.9
Basement	348	106.4	130.5	12.5	39.1	70.9	145.5	1762.5
<i>Type of building</i>								
Single	323	95.4	130.3	13.8	35.4	64.7	130.0	1762.5
Not single	577	130.5	137.9	12.5	46.0	84.2	187.8	1336.2
<i>Year construction/last renovation</i>								
Before 1990	525	117.0	134.6	12.5	41.0	74.4	161.4	1336.2
After 1990	375	119.2	138.6	16.0	42.6	79.0	160.9	1762.5
<i>Floor material</i>								
Other material	853	118.9	137.3	12.5	41.9	76.5	161.2	1762.5
Marble-granite	47	100.8	100.8	19.9	37.0	63.6	147.0	742.9

Fig. 3 Geo-lithological classification of the regional territory



$$\sum_{i=1}^n \rho_q\{r_{iq}\}, \tag{2}$$

288 Here, $r_{iq} = (y_i - \mathbf{x}_i^T \boldsymbol{\beta}_q) / \sigma$, σ is a scale parameter, the
 289 asymmetric loss function is $\rho_q\{r_{iq}\} = 2\rho\{r_{iq}\}$
 290 $[qI(r_{iq} > 0) + (1 - q)I(r_{iq} \leq 0)]$ and $I(\cdot)$ is the indicator
 291 function. The regression parameters could differ for dif-
 292 ferent values of q . M-quantile, quantile and expectile
 293 regression models can be obtained as special cases by using
 294 different specifications for the asymmetric loss function ρ .

295 See details in Bianchi et al. (2018). Throughout the paper
 296 the Huber loss function (Huber 1981) is used to define the
 297 linear M-quantile regression model:

$$\rho_q\{r_{iq}\} = 2 \begin{cases} (c|r_{iq}| - c^2/2)|q - I(r_{iq} \leq 0)| & |r_{iq}| > c \\ (r_{iq}^2/2)|q - I(r_{iq} \leq 0)| & |r_{iq}| \leq c, \end{cases} \tag{3}$$

299 where c is a tuning constant. Conventionally, in an M-reg-
 300 regression, the data analyst tunes this constant to provide a
 301 trade-off between robustness and efficiency. Huber (1981)

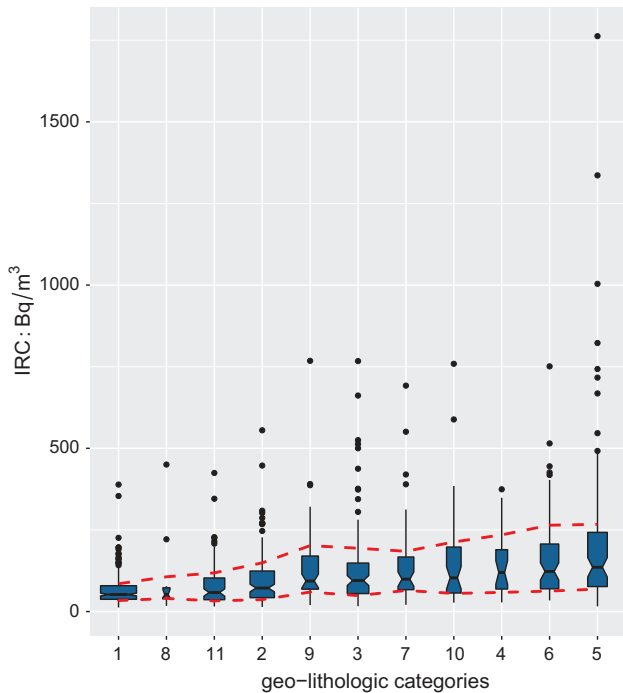


Fig. 4 Geo-lithological classification of the regional territory. Class labels: 1 Alluvial plain, 2 Foothill deposit, 3 Limestone, 4 Alluvial fan, 5 Debris, 6 Dolomite rocks, 7 Acid rocks, 8 Basic rocks, 9 Metamorphic rocks, 10 Alluvial plain, and 11 mountain valley. Dashed lines connect the 20th and 80th IRC quantiles of each geo-lithologic category. Larger boxes indicate a large sample size in the corresponding geo-lithologic class

302 proposes a value between 1 and 2. In particular, the author
 303 suggests 1.5. In the `r1m` function of the package `MASS` in R
 304 the default value for the tuning constant is 1.345. It cor-
 305 responds to 95% of the efficiency of the estimates under
 306 normality. It means that when the errors follow the normal
 307 distribution, setting c equal to a large value, say 100, is the
 308 most appropriate choice. In this case if a smaller value, say
 309 1.345, is used it will reduce the efficiency of the estimates
 310 because the tuning constant will offer unnecessary
 311 robustness. The tuning constant $c = 1.345$ is used
 312 throughout the paper. To overcome this ad hoc approach to
 313 selecting the tuning constant, Bianchi et al. (2018) exten-
 314 ded the data-driven method by Wang et al. (2007) to an
 315 M-quantile regression to estimate c via likelihood equa-
 316 tions. This approach could also be applied for the semi-
 317 parametric M-quantile models we propose in this paper, but
 318 this is beyond the scope of this work, however, and is left
 319 for future research.

320 Recently, M-quantile regression models that consider
 321 the two-level hierarchical structure in the data by including
 322 random effects were proposed by Tzavidis et al. (2016).
 323 The maximum likelihood method has been used by the
 324 authors for the estimation of model parameters.

The model suggested in this paper also includes a ran- 325
 326 dom intercept term to account for the clustering of data
 327 sharing the same geological substratum. However, as
 328 mentioned above, radon dynamics shows regularities when
 329 monitored in space both on the local or on the large-scale
 330 that are typically far from being linear. Hence, to adjust for
 331 these potential non-linear effects in the regression, we also
 332 add a flexible component to the linear predictor of the
 333 mixed effect M-quantile model. In particular, we use
 334 penalised splines. Penalised splines are effective tools for a
 335 number of reasons. Firstly, they are reasonably simple to
 336 implement, being a relatively straightforward extension of
 337 a linear M-quantile regression. Secondly, their flexibility
 338 enables the inclusion in a wide range of modelling features.
 339 More specifically, penalised splines account for spatial
 340 dependencies in the IRC data in the semiparametric
 341 M-quantile regression model adopted in the case study
 342 presented below. This component is expected to grasp not
 343 only large-scale dependencies of the radon data but also
 344 spatial local effects on the concentration field. Splines rely
 345 on a set of basis functions to handle non-linear structures in
 346 the data. In this paper, we assume that the spatial pattern of
 347 the variable of interest can be explained as a function of the
 348 location of a point that is represented by its cartographic
 349 coordinates. Thus, a bivariate smoothing spline is included
 350 in the additive specification of the model and is specified in
 351 terms of a set of bivariate basis functions. Following
 352 Ruppert et al. (2003), Pratesi et al. (2009) suggested the
 353 use of radial basis functions to derive low-rank thin plate
 354 splines.

The model we propose for a specified M-quantile q is: 355

$$MQ_y(q|\mathbf{X}, \mathbf{Z}, \mathbf{Z}_{sp}; \psi) = \mathbf{X}\boldsymbol{\beta}_q + \mathbf{Z}\mathbf{u}_q + \mathbf{Z}_{sp}\boldsymbol{\gamma}_q, \quad (4)$$

where \mathbf{X} is a matrix of dimension $n \times p$ of auxiliary vari- 357
 358 ables, $\boldsymbol{\beta}_q$ is the $p \times 1$ vector of M-quantile regression
 359 coefficients; \mathbf{u}_q is a $G \times 1$ vector of geological categories
 360 and $\boldsymbol{\gamma}_q$ is a $K \times 1$ vector of random effects associated with
 361 the spline matrix; \mathbf{Z} is an incidence $n \times G$ matrix coding
 362 the point-geological class hierarchy; \mathbf{Z}_{sp} is a $n \times K$ spline
 363 matrix and K is the number of spline knots. More specifi-
 364 cally (Opsomer et al. 2008),

$$\mathbf{Z}_{sp} = [C(\mathbf{w}_i - \mathbf{k}_j)]_{1 \leq i \leq n}^{-1/2} [C(\mathbf{k}_j - \mathbf{k}_k)]_{1 \leq j, k \leq K} \quad (5)$$

\mathbf{k}_j and \mathbf{k}_i , $j = 1, \dots, K$, $i = 1, \dots, K$, being two-dimen- 366
 367 sional vectors representing the cartographic coordinates of
 368 knots j and k . \mathbf{w}_i is a two-dimensional vector representing
 369 the cartographic coordinates of the sampling location i and
 370 $C(\mathbf{s}) = \|\mathbf{s}\|_2^2 \log \|\mathbf{s}\|_2$ where $\mathbf{s} \in \mathbb{R}^2$ and $\|\mathbf{s}\|_2$ is the Eucli-
 371 dean norm of \mathbf{s} in \mathbb{R}^2 .

Differently from the model suggested by Pratesi et al. 372
 373 (2009), we assume that the coefficients of the spline matrix
 374 in the linear predictor are random coefficients. A practical

375 advantage of the mixed model representation of the spline
 376 lies in fitting the model. The usual penalised spline-fitting
 377 criterion requires estimating a penalising or smoothing
 378 parameter prior to model estimation. Cross-validation is
 379 usually suggested as an appropriate way to tackle the
 380 problem. The mixed model representation avoids this step
 381 since the model can be estimated directly using routines
 382 that are appropriate for linear mixed models. Furthermore,
 383 including random coefficients for the spline basis compo-
 384 nents permits us to account for the bias due to omitted
 385 variables or unmeasured confounders. In addition, as
 386 advocated by Ruppert et al. (2003), treating the coefficients
 387 of the knots as random leads to a smoother representation
 388 of the estimated effect, compared to using fixed effects
 389 specification, and avoids data overfitting. As mentioned
 390 above, the spline component of the model is expected to
 391 catch the spatial regularity of IRC data and is used to
 392 visualise the results by smoothed maps. Hence, the random
 393 effect specification of the spline seems to be more appro-
 394 priate for this end. We define the following modified
 395 estimating equations, extending the idea of asymmetric
 396 weighting of the residuals to estimate the regression coef-
 397 ficients and the variance components (Tzavidis et al. 2016;
 398 Borgoni et al. 2018):

$$\mathbf{X}^T \mathbf{V}_q^{-1} \mathbf{U}_q^{1/2} \psi_q \{ \mathbf{r}_q \} = \mathbf{0} \tag{6}$$

$$\begin{aligned} & \frac{1}{2} \psi_q \{ \mathbf{r}_q \}^T \mathbf{U}_q^{1/2} \mathbf{V}_q^{-1} \mathbf{Z} \mathbf{Z}^T \mathbf{V}_q^{-1} \mathbf{U}_q^{1/2} \psi_q \{ \mathbf{r}_q \} \\ & \quad - \frac{K_{2q}}{2} \text{tr} \left[\mathbf{V}_q^{-1} \mathbf{Z} \mathbf{Z}^T \right] = 0 \\ & \frac{1}{2} \psi_q \{ \mathbf{r}_q \}^T \mathbf{U}_q^{1/2} \mathbf{V}_q^{-1} \mathbf{Z}_{sp} \mathbf{Z}_{sp}^T \mathbf{V}_q^{-1} \mathbf{U}_q^{1/2} \psi_q \{ \mathbf{r}_q \} \\ & \quad - \frac{K_{2q}}{2} \text{tr} \left[\mathbf{V}_q^{-1} \mathbf{Z}_{sp} \mathbf{Z}_{sp}^T \right] = 0 \\ & \frac{1}{2} \psi_q \{ \mathbf{r}_q \}^T \mathbf{U}_q^{1/2} \mathbf{V}_q^{-1} \mathbf{V}_q^{-1} \mathbf{U}_q^{1/2} \psi_q \{ \mathbf{r}_q \} - \frac{K_{2q}}{2} \text{tr} \left[\mathbf{V}_q^{-1} \right] = 0. \end{aligned} \tag{7}$$

402 Let $\mathbf{r}_q = \mathbf{U}_q^{-1/2} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}_q)$ denote the vector of scaled
 403 residuals with components r_{ijq} , \mathbf{U}_q the diagonal matrix with
 404 diagonal elements u_{ijq} equal to the diagonal elements of the
 405 covariance matrix \mathbf{V}_q and $\psi_q(r)$ the derivative of a loss
 406 function ρ_q . The covariance matrix \mathbf{V}_q is defined by
 407 $\mathbf{V}_q = \boldsymbol{\Sigma}_{\epsilon_q} + \mathbf{Z} \boldsymbol{\Sigma}_{u_q} \mathbf{Z}^T + \mathbf{Z}_{sp} \boldsymbol{\Sigma}_{\gamma_q} \mathbf{Z}_{sp}^T$, with $\boldsymbol{\Sigma}_{u_q} = \sigma_{u_q}^2 \mathbf{I}_G$,
 408 $\boldsymbol{\Sigma}_{\gamma_q} = \sigma_{\gamma_q}^2 \mathbf{I}_K$, and $\boldsymbol{\Sigma}_{\epsilon_q} = \sigma_{\epsilon_q}^2 \mathbf{I}_n$, where $\sigma_{u_q}^2$, $\sigma_{\gamma_q}^2$ and $\sigma_{\epsilon_q}^2$ are the
 409 quantile-specific variance components. \mathbf{I}_n is an identity
 410 matrix of size n and $K_{2q} = E[\psi_q(\boldsymbol{\epsilon}) \psi_q(\boldsymbol{\epsilon})^T]$ with
 411 $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I}_n)$. To obtain estimators of $\boldsymbol{\beta}_q$, $\sigma_{u_q}^2$, $\sigma_{\gamma_q}^2$, $\sigma_{\epsilon_q}^2$,
 412 Eqs. (6) and (7) are solved iteratively. For Eq. (6) a
 413 Newton–Raphson algorithm is used and for (7) the fixed-
 414 point iterative method is implemented to get the estimates.
 415 The algorithm is implemented by the authors in a function

in the R software (R Core Team 2017). A sandwich esti-
 mator is adopted to make inference on the model param-
 eters. Details of the estimation algorithm and the variance
 estimators are reported in Tzavidis et al. (2016).

4 M-quantile modelling of geocoded radon data

In this section, the model discussed in Sect. 3 is applied to
 the IRC data presented in Sect. 2. The set of covariates
 introduced in Table 1 that may potentially have an effect
 on IRC are those included in the M-quantile models con-
 sidered hereafter.

As mentioned in the introduction of this paper, IRC
 tends to vary across space showing regular possibly non-
 linear patterns of values due to a number of environmental,
 geological and anthropic factors. Geographical coordinates
 of the measurement locations can be considered as a sur-
 rogate variable of all these factors that may happen to be
 unmeasured or even unmeasurable. In this paper, we pro-
 pose including this component in three ways. Firstly, the
 IRC of the study region tends to be quite different for the
 south compared to the north (Borgoni et al. 2011), hence, a
 trend surface model (Cade et al. 2005; Koener and Mizera
 2004) is specified semiparametrically by a bivariate thin
 plate spline transformation of the cartographic coordinates,
 as discussed in Sect. 3. Secondly, IRC values detected in
 buildings that are built on the same type of soil can be
 expected to be more similar than those detected in different
 geo-lithological areas. As the data show a hierarchical
 structure, we explicitly consider this aspect in the model
 specification, including a random effect to capture the
 variability within geological areas, ending up in the semi-
 parametric M-quantile random effect model presented in
 Sect. 3. Thirdly, since, as mentioned above, high IRC can
 also be found in areas where faults are present, the distance
 to the nearest tectonic fragment is also included in the
 model.

In this section, we investigate two different issues
 related to IRC modelling: (1) the identification of radon-
 prone areas, and (2) the identification of those character-
 istics that make a building more exposed to high IRC. The
 latter analysis also allows for an identification of those
 building profiles that are more exposed to higher concen-
 tration and for which it can be appropriate to plan reme-
 diation actions or actions that are able to prevent gas
 accumulation.

Some preliminary analyses reported in Appendix 1 of
 this paper clearly show that the data may contain outliers
 that prevent the Gaussian assumptions from being met,
 leading to biased and inefficient estimates of the model
 parameters. Several papers (see Huggins 1993; Huggins

466 and Loesch 1998) suggested the robust estimation of mixed
467 models to protect against departures from normality. This
468 can be obtained using a loss function in the log-likelihood
469 that increases with the regression residuals at a slower rate
470 than the squared loss function. As described in the previous
471 sections, resorting to the M-quantile approach permits us to
472 robustly estimate the relationship between IRC and a set of
473 covariates in a natural manner.

474 4.1 Radon-prone areas: identification 475 and mapping

476 Health radon-related concerns have generated a growing
477 interest in identifying those regions of the territory where
478 high IRC are expected, the so-called ‘radon-prone areas’
479 (RPA). A number of different approaches have been sug-
480 gested for this, mostly based on various cluster detection
481 methods adopted to either delineate spatial clusters or
482 improve the understanding of the spatial dynamic of radon
483 using an automatic detection of those regions of space that
484 are ‘anomalous’, ‘unexpected’ or otherwise ‘interesting’
485 (Sarraf et al. 2016).

486 According to the World Health Organization (2009),
487 various definitions of radon-prone areas exist. Typically,
488 countries define RPA as regions where the estimated per-
489 centage of homes, whose radon concentrations exceed a
490 reference value τ , oversteps a threshold q . The Italian
491 legislation is compliant to the WHO suggestion to define
492 RPA as regions where ‘there is a high probability of finding
493 high (indoor) radon concentrations’ (art. 10-ter, comma 2,
494 D.L.vo 241/00). The above-mentioned definitions suggest
495 that an RPA is a region where $P(\text{IRC} > \tau)$ is high: high
496 meaning $P(\text{IRC} > \tau) > q$ and with q denoting a fixed
497 threshold, although different reference levels are suggested
498 by different local authorities. Below, we exemplify the
499 procedure using $q = 0.15$. Indicating by ξ_{1-q} the $(1 - q)$ -
500 quantile of IRC, this also means that $\xi_{1-q} > \tau$. Hence, one
501 can equivalently define an RPA as a region where a suf-
502 ficiently high-order quantile of IRC is above the reference
503 level τ .

504 Moving from this definition, RPA identification based
505 on conditional quantiles of the radon distribution sounds
506 more appropriate than basing it on cluster algorithms.
507 Borgoni et al. (2010) also suggested a quantile-based
508 approach adopting conventional kriging procedures cou-
509 pled with Monte Carlo sequential (Gaussian) simulations to
510 approximate the conditional distribution of radon at each
511 point in space. Directly modelling the tail of the distribu-
512 tion using an M-quantile or a quantile approach (Fontanella
513 et al. 2015), seems, however, a more direct and natural way
514 to operate.

515 For the rest of this section we considered the 85th
516 M-quantile as the reference level to identify RPA. In order
517 to estimate such an M-quantile at different locations in
518 space, the following semiparametric M-quantile random
519 effect regression model is employed:

$$MQ_{y_i}(0.85|d_{ij}, x_{1ij}, x_{2ij}, \mathbf{z}_i, \mathbf{z}_{spi}; \psi) = \alpha + d_{ij}\delta + x_{1ij}\beta_1 + x_{2ij}\beta_2 + \mathbf{z}_i\mathbf{u} + \mathbf{z}_{spi}\gamma, \quad (8)$$

where

- d is the fault distance;
- x_1 and x_2 are the coordinates (respectively, longitude and latitude) of the measurement points in UTM projection;
- \mathbf{z}_i is a vector of geo-lithologic class indicators, i.e 0–1 variables;
- \mathbf{z}_{spi} is the row of the Z_{sp} matrix defined in Eq. 5 pertinent to sampling dwelling i located in geo-lithologic class j .

521 We use 50 knots for the spline, obtained by applying the
522 partitioning clustering algorithm CLARA (Kaufman and
523 Rousseeuw 1990) to the sampling locations. The estimated
524 parameters are reported in Table 2. The linear effect of the
525 cartographic coordinates has been found to be not signifi-
526 cant. However, the estimate of the variance of γ is about
527 three times its estimated standard error pointing out a
528 strong effect of the random component of the spline. This
529 demonstrates that the spatial variability due to location is
530 definitely relevant in the IRC dynamic. We notice that the
531 two linear terms are retained in the model in the following
532 analysis to correctly specify the spline component.

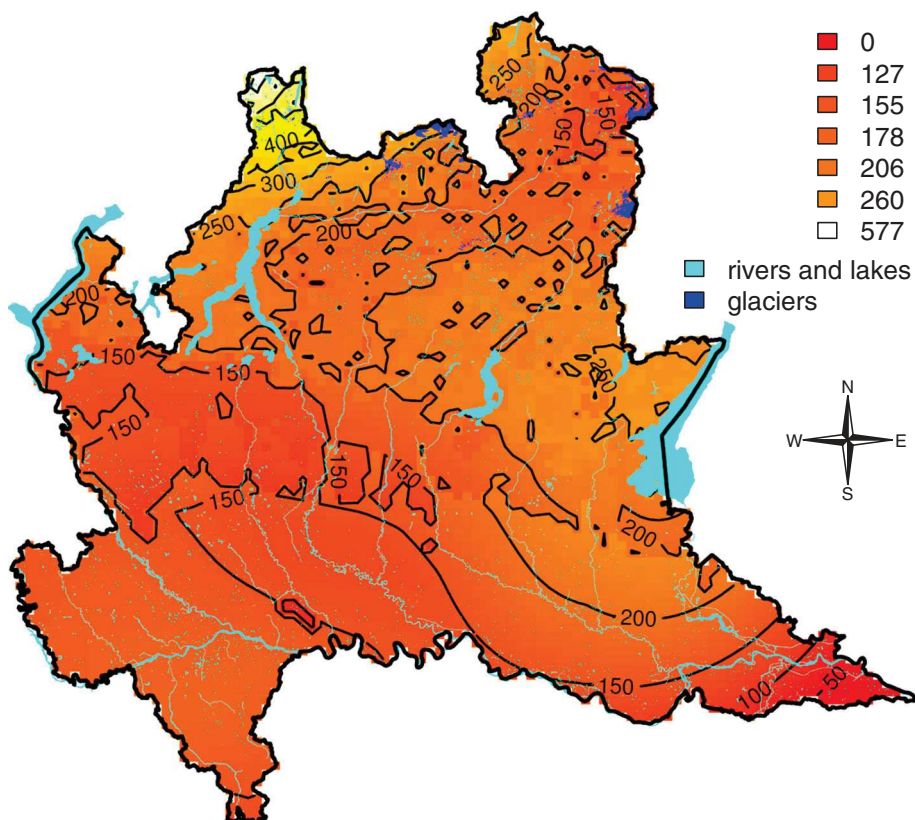
533 Since the aim of this analysis is to classify the areas of
534 the region according to their proneness to radon without
535 considering any particular building typology, the structural
536 and architectonic characteristics of the building are not
537 included in the model. The issue of assessing the impact of
538 building factors will be addressed in the following section.

539 The estimates of the 85th M-quantile are depicted in
540 Fig. 5 where the surface has been discretised via a grid of

Table 2 Estimates of the RPA model

	Estimate	Std. error	p value
Intercept (α)	178.75	275.62	0.52
Fault distance (δ)	− 76.14	44.01	0.08
Longitude (β_1)	237.45	385.68	0.54
Latitude (β_2)	− 154.31	387.98	0.69
σ_e^2 (individual)	7758.4	220.07	
σ_u^2 (geo-lithology)	302.4	402.06	
σ_γ^2 (spline)	75250.2	28397.95	

Fig. 5 Fitted surface of the 85th M-quantile. The values in the legend represent the minimum, the maximum and five equispaced percentiles of the 85th M-quantiles estimated by the model in Eq. (8) at the grid points



551 4651 points internal to the administrative boundary of the
 552 study region. For each point of the grid, the geo-lithologic
 553 class has been retrieved by overlaying the point on the map
 554 in Fig. 3 and the distance to its nearest tectonic lineament
 555 has been calculated by applying Eq. (1). Analogously, the
 556 spline base z_{sp} has been calculated for each location of the
 557 grid by applying Eq. (5). As far as the fixed effects are
 558 concerned, the value of the covariates are multiplied by the
 559 estimated coefficients shown in Table 2, whereas the ran-
 560 dom-effects vector \mathbf{u} and γ in Eq. (8) have been predicted
 561 using a modified Fellner equation (Fellner 1986) as pro-
 562 posed by Borgoni et al. (2018) for the three-level
 563 M-quantile random effect models. The issue of predicting
 564 random effects is also discussed by Geraci and Bottai
 565 (2014) and Tzavidis et al. (2016). The estimated M-quan-
 566 tiles are calculated by summing the different components
 567 according to Eq. (8) and a raster of 4,651 pixels are
 568 eventually obtained and are displayed in Fig. 5.

569 In order to identify radon-prone areas, Fig. 5 has been
 570 transformed into a binary map by colouring those pixels in
 571 red, where the estimated M-quantile is above the reference
 572 level. Figure 6a, b show the results using a reference level
 573 τ corresponding to 200 Bq/m^3 and 300 Bq/m^3 , respec-
 574 tively. The latter reference level has been suggested by the
 575 recent 2013/59/EURATOM European recommendation as
 576 a suitable reference value for the annual average indoor

concentration of radon, whereas the former was suggested
 by the 90/143/Euratom recommendation and it was widely
 used in the past.

4.2 Assessing the role of influential factors on IRC

As noticed in Sect. 2 there are a number of factors in
 addition to space and geological dimensions that can
 potentially affect the concentration of radon in an indoor
 environment, such as building-specific characteristics. The
 exploratory analysis also suggests that the impact of a
 given characteristic can be different at different concen-
 tration levels. Hereafter, the conditional distribution of IRC
 is modelled using the approach introduced in Sect. 3 as a
 function of these building factors in order to quantify their
 potential effects and how they differ at different levels of
 IRC. Hence, the model in Eq. 8 is expanded to include
 these covariates as fixed effects. The baseline house is
 located in a building in direct contact with the ground,
 equipped with an air conditioning system and constructed
 or refurbished in 1990 or before with walls made by
 materials other than stone and floors made by materials
 other than marble or granite.

Table 3 shows the estimated parameters for three dif-
 ferent M-quantiles, 0.25, 0.5 and 0.75. As in the section

Author Proof

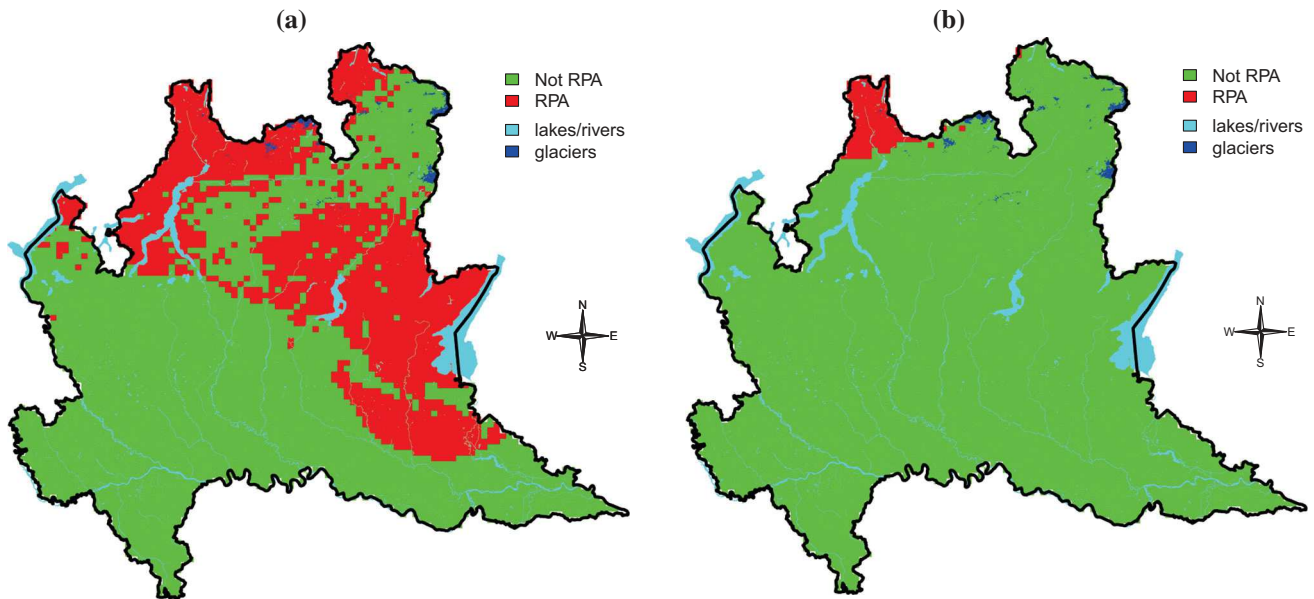


Fig. 6 Radon-prone areas (RPA) according to a reference value of a 200 Bq/m^3 and b 300 Bq/m^3

Table 3 Results—Semiparametric M-quantile random effect model: point estimates with standard errors in parentheses

	$q = 0.25$		$q = 0.50$		$q = 0.75$	
	Estimate	p value	Estimate	p value	Estimate	p value
Intercept	27.24 (13.53)	0.044	45.59 (25.12)	0.069	84.93 (74.86)	0.256
Distance to nearest fault	- 3.21 (8.27)	0.697	- 10.25 (13.47)	0.446	- 35.19 (28.08)	0.21
Floor: marble or granite	- 5.41 (5.54)	0.328	- 8.88 (8.87)	0.316	- 22.43 (17.32)	0.195
Wall: stone	3.77 (3.94)	0.338	6.53 (6.31)	0.301	16.86 (12.38)	0.173
Years of construction/last renovation: after 1990 single buildings	7.99 (2.55)	0.001	11.25 (4.09)	0.005	14.11 (8.003)	0.077
Not in contact with the ground no air conditioning	6.35 (2.69)	0.018	10.36 (4.32)	0.016	24.98 (8.49)	0.003
	- 5.73 (2.67)	0.032	- 9.74 (4.28)	0.022	- 18.24 (8.35)	0.029
Longitude	1.09 (5.16)	0.832	- 1.48 (8.29)	0.858	- 7.59 (16.24)	0.64
	33.4 (15.44)	0.03	42.43 (30.23)	0.16	87.73 (97.85)	0.369
Latitude	36.62 (18.77)	0.05	48.81 (35.19)	0.165	38.58 (106.14)	0.716
	σ_ϵ^2 (individual)	819.88 (48.91)		3172.99 (241.93)		6283.1 (421.41)
σ_u^2 (geo-lithology)	42.21 (26.41)		149.7 (85.48)		220.13 (92.08)	
σ_γ^2 (spline)	75.96		491.78		3140.52	

601 above, 50 knots have been used for the radial spline,
602 obtained by applying the CLARA algorithm to the sam-
603 pling locations. Appendix 2 provides a short sensitivity
604 analysis where the estimated parameters of quantile and
605 M-quantile regression models are compared. Although the
606 two approaches cannot be directly compared, as these
607 models target different location parameters, the results
608 show that the coefficients of the M-quantile regression
609 model are in the same direction as the ones based on
610 quantile regressions.

611 Figure 7 shows the estimated effect by M-quantile of
612 each covariate that we have included in the model. Con-
613 fidence bands across the M-quantiles are also reported in
614 the graphs to display the sampling variation. For each
615 considered M-quantile, the band is obtained by calculating
616 the point-wise 95% confidence interval of the regression
617 coefficients and it is displayed in the graph by a grey-
618 shaded area around the line. It can be seen that the variation
619 between M-quantiles is diverse, sometimes even exten-
620 sively so, and it tends to increase at the edges of the
621 M-quantile order, where such an increase can be quite
622 large. This is quite typical for quantile modelling, since
623 estimates too far from the centre of the distribution usually
624 cannot be determined with high precision. This problem
625 can possibly be exacerbated by the hierarchical structure of
626 the data. The clustering of the measurement points (about
627 900 in all) in 11 classes implies that the tail of the distri-
628 butions cannot be well frequented by the data, as is also
629 shown by the conditional boxplot in Fig. 4, contributing to
630 reducing the information.

631 Concerning the random components of the model at the
632 lower quantile, a large part of the variation is due to
633 individual variability whereas both the variability due to
634 geo-lithology and space are definitely minor, as one could
635 expect. Moving towards higher quantiles, the spatial
636 component tends to become more and more relevant.
637 Figure 8 shows the estimated spline effects at the three
638 considered M-quantiles. The estimated effects are obvi-
639 ously larger at higher quantile orders. The maps also show
640 that a substantial homogeneity in space exists at the lower
641 quartile (Fig. 8a) apart from some picks in the mountains in
642 the far north and south-east of the region. At high orders,
643 the spatial dynamics tend to vary more due to the large-
644 scale tendency over the investigated region and due to local
645 effects caught by the semiparametric component of the
646 model.

647 We finally observe that an important part of spatial
648 prediction refers to the measurement of uncertainty. The
649 predicted spline effects at M-quantiles have an uncertainty
650 that could be estimated following Ruppert et al. (2003) and
651 Opsomer et al. (2008). As Ruppert et al. (2003) noticed,
652 the mixed model formulation of penalised splines is a
653 convenient artefact for estimating the smoothing

parameters while the ML or REML variance component 654
estimation provides estimates of the smoothing parameter 655
that generally behaves quite well. The standard errors 656
derived according to the approach suggested by these 657
authors are expected to account for both the error compo- 658
nents (variance and squared bias) and can be somewhat 659
wider than those obtained without using a mixed model 660
representation. A similar approach can also be adopted to 661
calculate the standard error for the predicted bivariate 662
spline effects at M-quantiles. However, this is beyond the 663
scope of this work and it is left for further research. 664

665 Quite surprisingly, the geo-lithological component
666 remains negligible even at the highest quantile. This can be
667 due to the spatial resolution of the geological and litho-
668 logical information. The maps available for this analysis
669 are scaled 1:250,000, hence, different geological structures
670 can be mixed up in different classes because of a low
671 resolution inducing the inhomogeneity of the geological
672 units, and the geo-lithological effect can be watered down.

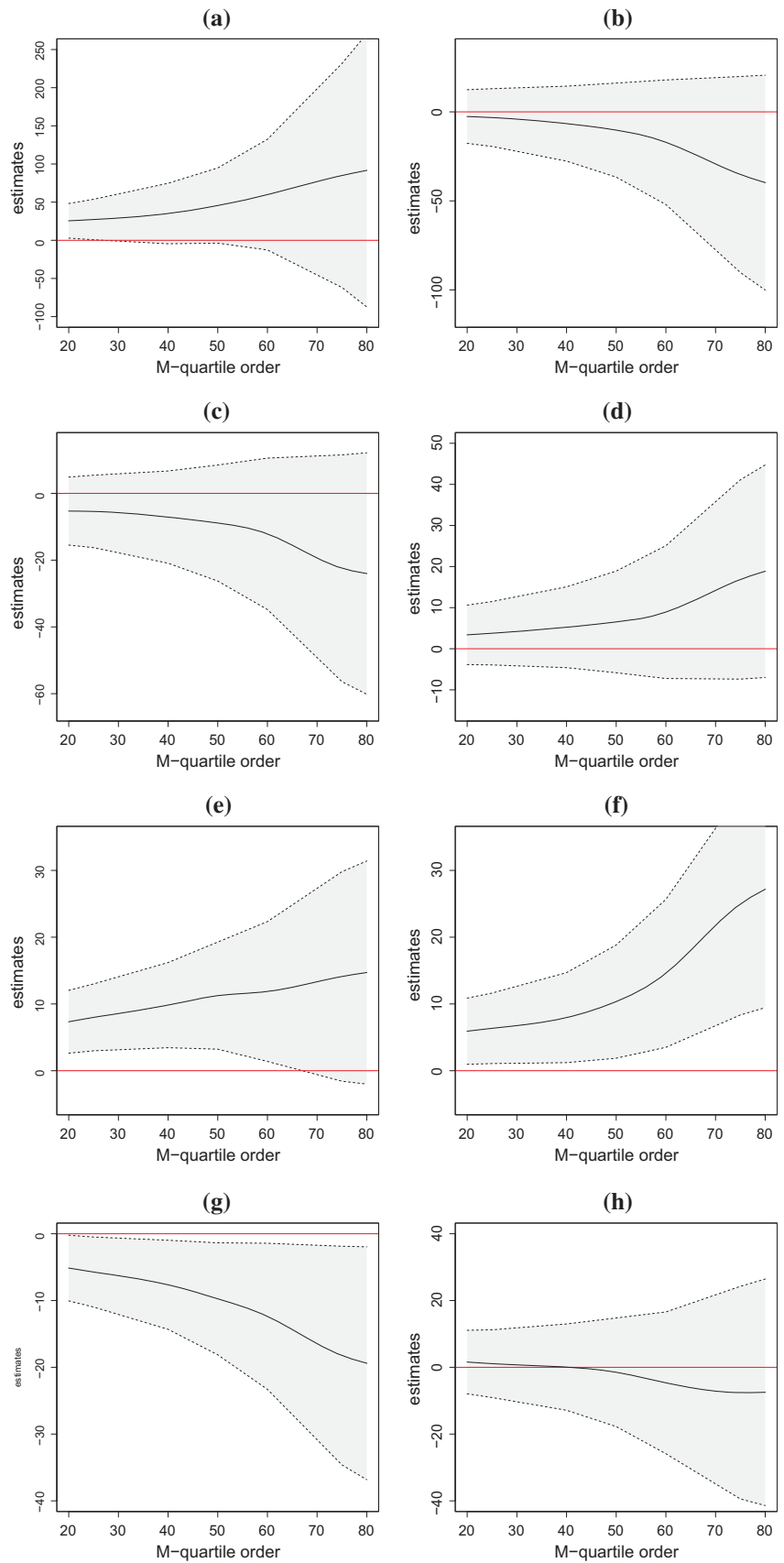
673 5 Discussion and conclusions

674 In this paper, a semiparametric random-effect M-quantile
675 model is introduced in order to investigate radon concen-
676 tration within buildings. It has been shown how the model
677 can be estimated within the framework of robust maximum
678 likelihood by using a numerical optimisation method based
679 on the Newton-Raphson and fixed-point algorithms that
680 apply the data of an indoor radon gas monitoring survey
681 carried out by the Agency of Environmental Protection of
682 Lombardy (Italy) in 2003.

683 The proposed model allows for an investigation of how
684 a set of covariates acts at different levels (M-quantiles)
685 of the IRC distribution while accounting for the hierarchi-
686 cal nature of the data. In particular, a set of building char-
687 acteristics are included in the model as well as information
688 concerning the geological structure of the soil.

689 One of the building characteristics that was found to be
690 statistically significant at M-quantiles was whether the
691 building is in direct contact with the ground and the
692 building type. Buildings in contact with the ground and
693 detached buildings are found to have a higher indoor
694 concentration than other buildings and the impact of those
695 variables becomes larger as the M-quantile order increases
696 (i.e. for those situations more seriously affected by large
697 concentrations). This is not an unexpected result. Unlike
698 many other indoor air pollutants that are correlated to
699 outdoor air pollution, radon gas concentrations in homes
700 are related primarily to the ingress of radon from ground
701 sources. Hence, being in contact with the ground fosters
702 gas accumulation. Condominiums are often constructed out
703 of concrete. The radium content of the concrete is typically

Fig. 7 Estimated coefficients of M-quantile regressions: **a** intercept, **b** fault distance, **c** floor material, **d** wall material, **e** year from construction/last renovation, **f** single building, **g** not in contact with the ground, **h** air conditioning system. Shaded areas represent the 95% confidence intervals



Author Proof

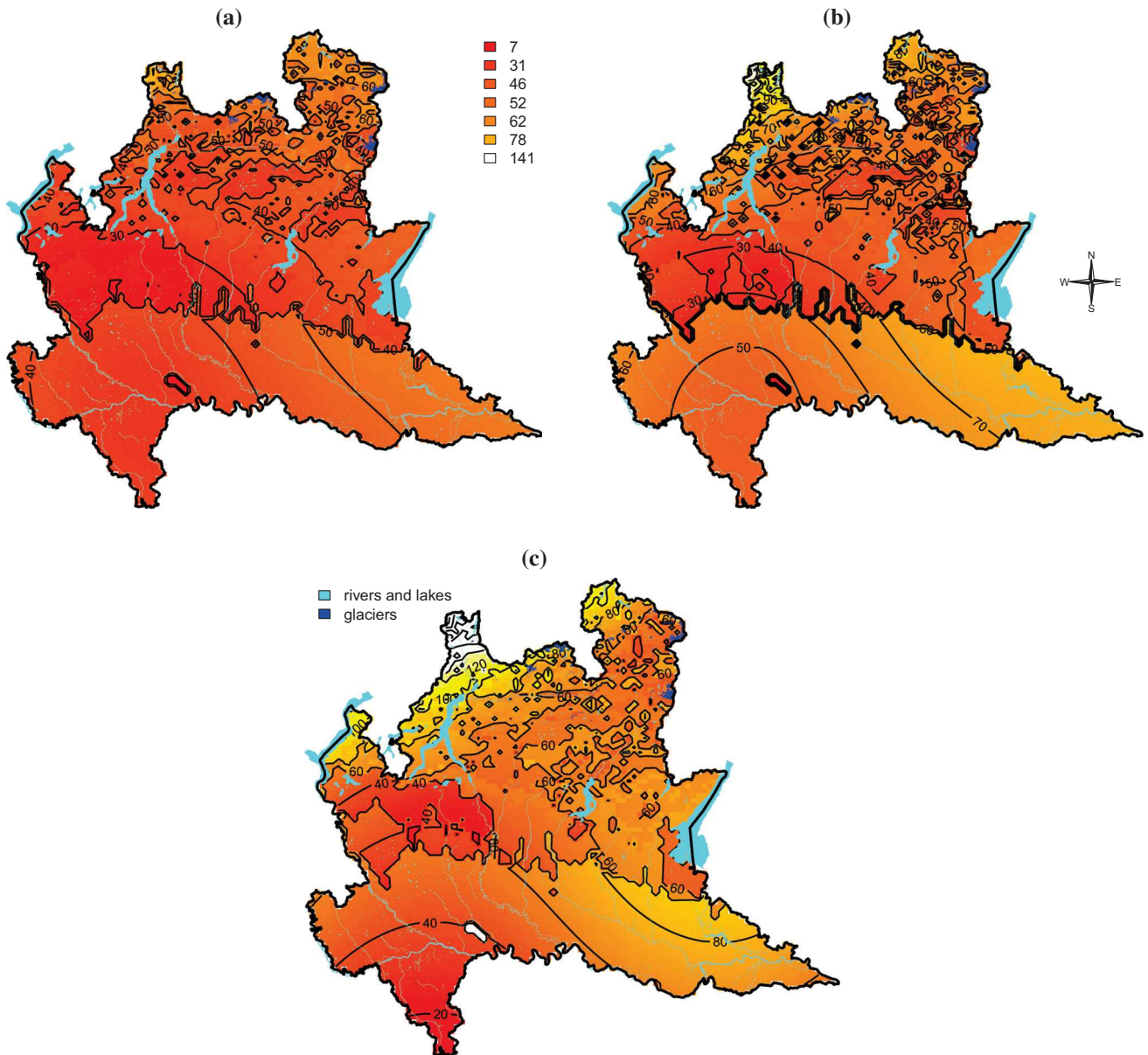


Fig. 8 Spline effects at the three considered M-quantiles: a 25th, b median, and c 75th

704 not high and this may explain why high-density housing is
 705 characterised by lower radon concentration. The statistical
 706 significance of the regression coefficient associated with
 707 whether the building was built or refurbished after 1990 is
 708 mild and not completely clear at all M-quantile orders.
 709 There is widespread belief that increased weatherproofing
 710 and the energy efficiency of homes significantly contribute
 711 to the increase in residential radon concentrations.
 712 Nonetheless, uncertainty remains about their actual impact,
 713 in particular, whether energy efficiency guidelines include
 714 the consideration of air exchange rates and ventilation.
 715 Finally, although not all the variables considered in the
 716 model have been found to be statistically significant, it can

be observed that the estimates tend to be larger in modulus
 moving towards higher M-quantiles, suggesting that
 building characteristics can be expected to be effective
 potential levers for moderating critical situations.

Our findings provide useful indications in this direction,
 helping to identify those factors that mainly foster high
 concentration levels of the pollutant on a large and inho-
 mogeneous territory with several different house typolo-
 gies. Using the estimated regression coefficients, it is
 possible to classify the different typologies of buildings
 based on a selected M-quantile of the IRC distribution and
 to provide a ranking of the dwellings according to their
 proneness to IRC. To this end, we considered the

717
 718
 719
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729

M-quantile regression of order $q = 0.75$ that models how dwelling characteristics impact at a high level of the pollutant distribution. More specifically, fixating the spatial location of a building to a prefixed value, we combined the variables found to be statistically significant at the M-quantile $q = 0.75$ (p values smaller than 0.05), namely whether the building is in contact with the ground and whether it is a single building, with the geo-lithological classes obtaining a set of $2 \times 2 \times 11 = 44$ different building profiles for which the considered M-quantile has been estimated. For instance, for a single building that is not in contact with the ground located in a debris class, replacing the unknown parameters in Eq. (4) with their estimates (see Table 3) and ignoring the spatial component of the model, gives:

$$\widehat{MQ}_y(0.75) = 84.93 + 24.98 \times I(\text{single building} = 1) + (-18.24) \times I(\text{not in contact with the ground} = 1) + 20.89 \times I(\text{geo-lithological class} = \text{debris}) = 112.55$$

where 20.89 is the estimated residual for the debris class. These profiles are ranked according to their estimated value of the 0.75 M-quantile from the highest to the lowest, obtaining a measure of the building's proneness to a high IRC level. The top-five building profiles most prone to IRC are listed in Table 4.

Selecting an arbitrary location to provide the different scenarios can be done without loss of generality. Although the actual location does affect the estimated IRC M-quantile, the ranking of profiles provided by the M-quantile model based on the building characteristics does not change, considering other locations in space given the additive nature of the model. Looking at the five profiles most prone to high IRC listed in Table 4 we found that all of them are single buildings and four out of five are in contact with the ground and located on porous soils or soils characterised by weathered carbonate rocks (such as debris and dolomite) where radon emanation is known to be high despite the low concentrations of uranium. Hence, our results can help local authorities involved in environmental protection both to identify some guidelines for new buildings and to identify those dwelling typologies already present in the territory that should be monitored in order to mitigate concentration levels.

Given the serious health-related problems induced by the exposure of humans to radon gas, the usage of monitoring surveys for the identification of areas more prone to high IRC, named radon-prone areas above, has been promoted in many countries worldwide. We demonstrated how the semiparametric M-quantile model proposed in this paper provides a natural way to identify such areas flexibly and effectively, taking into account (1) the spatial dynamic of IRC via flexible bivariate spline transformations, and (2) information related to the geological and geophysical information included in fixed and random components of the model. It is worth noticing that the information concerning the geology of the soil is conveyed via digital maps that can be linked to the dataset by GIS operations. Hence, the spatial resolution of this information may have an impact on the precision of the estimates, and having high-resolution maps for those dimensions may sensibly improve the outcome of the analysis.

We also show how the outcome of the model can be visualised by using thematic maps. Such maps inform people and local authorities of where higher concentrations can be expected. Local authorities can use the maps to differentiate construction requirements according to different locations.

Finally, we recognise that when considering a complex phenomenon such as radon gas accumulation, the set of relevant factors may be larger than the one considered in the present paper, and adding these further control variables might lead to improved results. In particular, as has been demonstrated, for instance by Kemski et al. (2009), soil radon measurements are often considered to be a potential predictor of indoor concentration since soil gas containing radon leaks into houses through cracks or holes in the foundations because of the lower air pressure observed indoors compared to outside. Radiometric data has sometimes also been used to account for the radioactivity of the soil.

Numerous weather-related factors influence the ingress of radon into buildings, including wind, barometric pressure, rainfall, and indoor and outdoor difference of temperature variations (Rowe et al. 2002). Increased wind can exert small pressure differences between the lower levels of a dwelling and the outdoors and an increased precipitation can act to impede radon emanation. Climate

Table 4 The top-five building profiles most prone to IRC

Ranking	Contact with the ground	Building type	Geo-lithological class
1	In contact	Single	Debris
2	In contact	Single	Dolomite rocks
3	In contact	Single	Alluvial fan
4	Not in contact	Single	Debris
5	In contact	Single	Limestone

814 parameters affect the ventilation of indoor environment as
 815 well, in turn influencing the indoor concentration of the
 816 pollutant. Unfortunately, we did not have this information
 817 at hand for this study, and hence their effect has not been
 818 investigated, although it can easily be added to the model
 819 proposed in this paper, if available.

820 **Acknowledgements** The work of Nicola Salvati has been developed
 821 under the support of the Progetto di Ricerca di Ateneo 'From survey-
 822 based to register-based statistics: a paradigm shift using latent vari-
 823 able models' (Grant PRA2018-9). The authors were further supported
 824 by the MIUR-DAAD Joint Mobility Program (57265468).

825 **Appendix A: Preliminary data analysis**

826 Hereafter some preliminary data analyses is reported that
 827 motivates the need for a robust approach when modelling
 828 IRC data. To this aim an ordinary random effect model for
 829 the mean IRC that reflects the hierarchical structure of the
 830 data with buildings nested in the geological classes has
 831 been fitted using the function `lmer` of the R package
 832 `lme4`. Figure 9a shows the normal qq-plot of the indi-
 833 vidual residuals (i.e. residuals pertinent to the building
 834 level) whereas Fig. 9b displays the normal qq-plot of the
 835 residuals estimated from the model at the geological class
 836 level. These plots show that the normality assumptions of
 837 the ordinary mixed model are violated, which is also
 838 confirmed by the Shapiro-Wilk test (p values=0.0000078
 839 for the geological class residuals and p value= $2.2e-16$ for
 840 the building residuals). Figure 10a shows the histogram of
 841 the standardised building residuals obtained by the random
 842 effect regression model, whereas Fig. 10b displays the
 843 distribution of the standardised residuals by geological

844 classes. The histogram appears very skewed and some
 845 classes have many large positive residuals (larger than 2).
 846 Thus, influential observations seem to be present in the
 847 data. This is also confirmed by Fig. 11 that displays the
 848 Cook's Distance for the two sets of residuals.

849 It is clear that the data may contain outliers and influ-
 850 ential points that invalidate the Gaussian assumptions. In
 851 these circumstances, estimates of the model parameters are
 852 biased and inefficient and the robust approach suggested in
 853 this paper sounds more appropriate.

854 **Appendix B: Additional results for modelling
 855 geocoded radon data**

856 Appendix 1 provides a short comparison of the estimated
 857 parameters obtained from quantile and M-quantile regres-
 858 sion models. The two approaches cannot be directly com-
 859 pared since they target different location parameters.
 860 However, both approaches try to model location paramet-
 861 ers that are related to the same part of the conditional
 862 distribution of IRC. Table 5 reports the estimated regres-
 863 sion coefficients for $q = 0.5$ for two approaches: (1) the
 864 proposed semiparametric M-quantile random effect
 865 regression model (semiMQRE), and (2) a semiparametric
 866 quantile regression model (semiQR). semiQR is based on
 867 an additive quantile regression model (Koenker et al. 1994)
 868 where the spatial structure is captured by bivariate splines
 869 but without accounting for the hierarchical structure in the
 870 data by a random component. The results indicate that the
 871 coefficients based on M-quantile regression models are in
 872 the same direction as the ones based on quantile regression.
 873 However, with quantile regression convergence problems

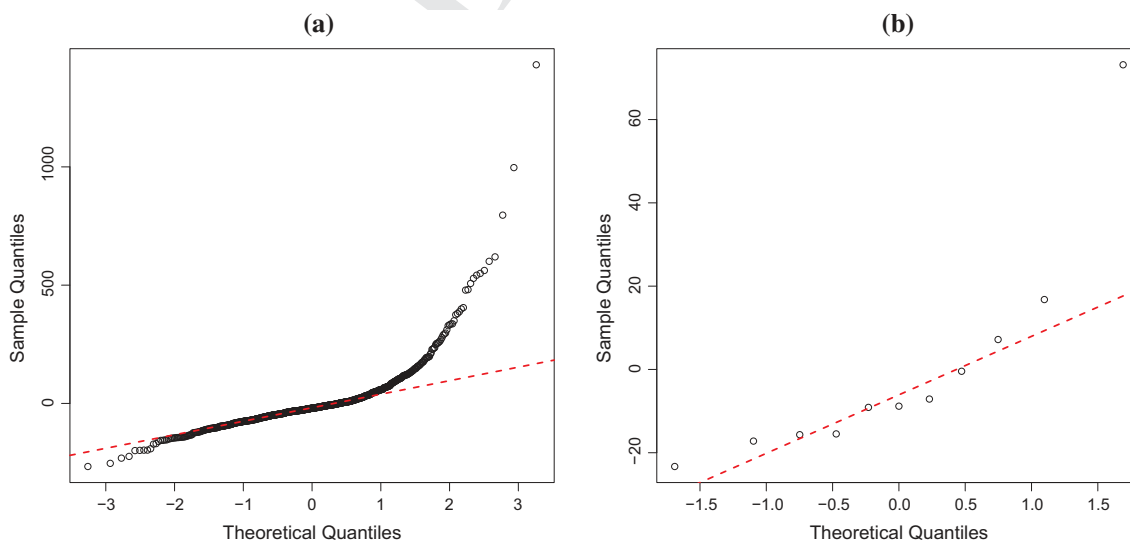


Fig. 9 QQ-plot of building residuals (a) and of geological class residuals (b) estimated by the two-level random effect regression model for the mean IRC

Author Proof

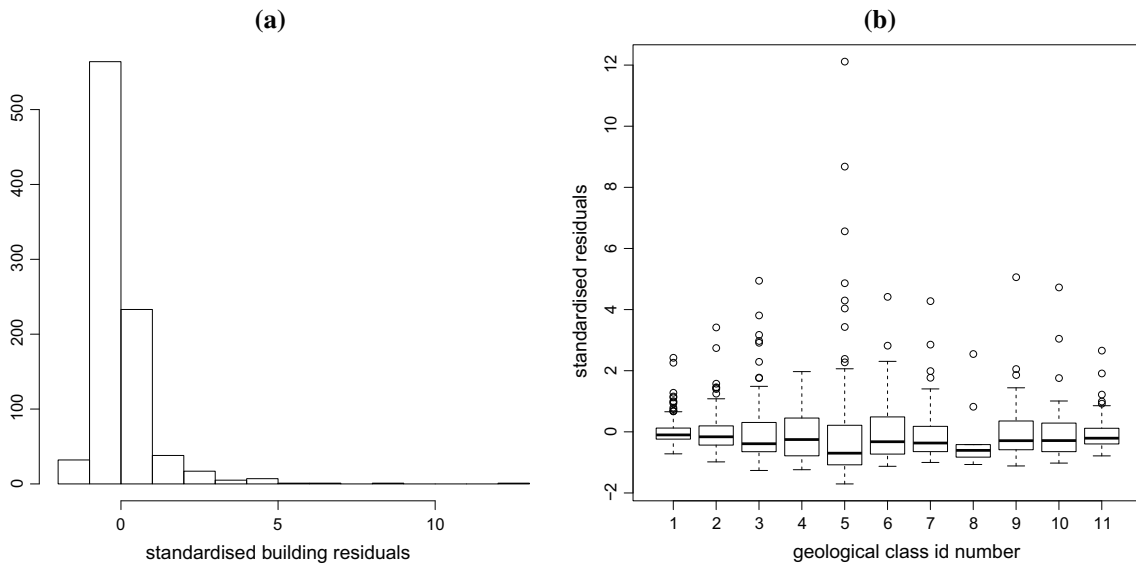


Fig. 10 Histogram of standardised building residuals of the two-level random effect regression model for the mean IRC (a); boxplots of standardised building residuals by geological classes (b)

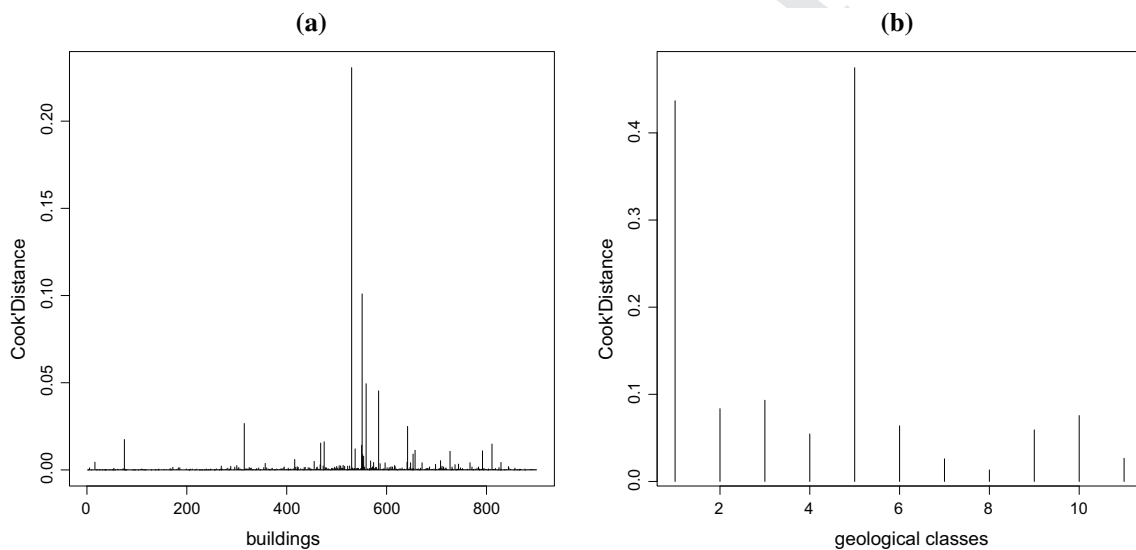


Fig. 11 Cook Distance of building residuals (a) and of geological class residuals (b) estimated by the two-level random effect regression model for the mean IRC

Table 5 Results—Semiparametric M-quantile and quantile regression models for $q = 0.5$: Point estimates with standard errors in parentheses

	semiMQRE		semiQR	
	Estimate	<i>p</i> value	Estimate	<i>p</i> value
Intercept	45.59 (25.12)	0.069	44.56 (11.02)	0.000
Distance to nearest fault	- 10.25 (13.47)	0.446	- 8.14 (13.71)	0.553
Floor: marble or granite	- 8.88 (8.87)	0.316	- 7.41 (8.90)	0.406

Table 5 (continued)

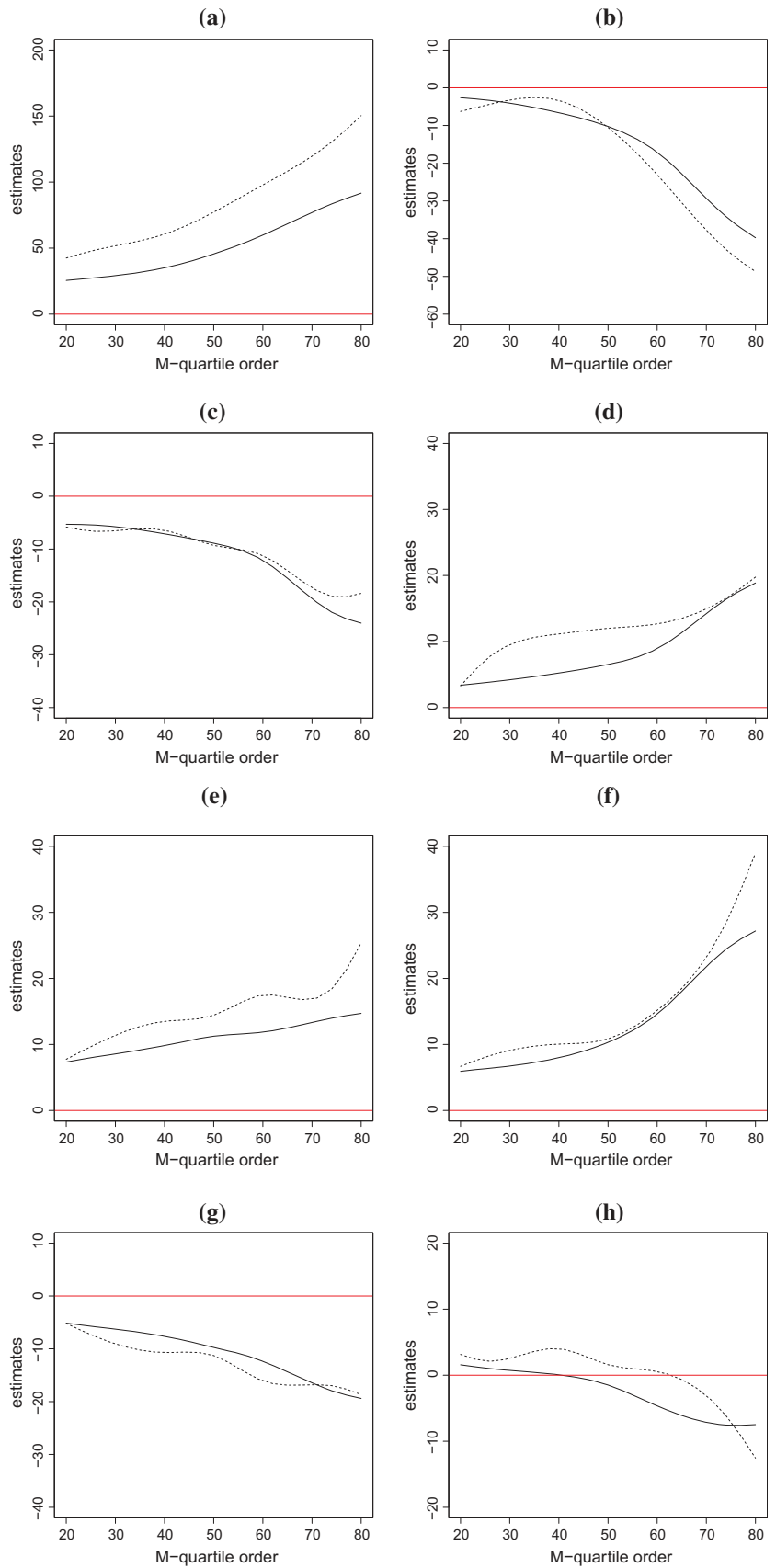
	semiMQRE		semiQR	
	Estimate	<i>p</i> value	Estimate	<i>p</i> value
Wall: stone	6.53 (6.31)	0.301	- 1.98 (9.07)	0.827
Years of construction/last renovation: after 1990 single buildings	11.25 (4.09)	0.005	7.61 (4.78)	0.112
	10.36 (4.32)	0.016	9.62 (4.74)	0.043
Not in contact with the ground no air conditioning	- 9.74 (4.28)	0.022	- 9.63 (4.87)	0.048
	- 1.48 (8.29)	0.858	- 1.10 (7.48)	0.883
Longitude	42.43 (30.23)	0.16	-	-
Latitude	48.81 (35.19)	0.165	-	-
σ_ϵ^2 (individual)	3172.99 (241.93)		-	
σ_u^2 (geo-lithology)	149.7 (85.48)		-	
σ_γ^2 (spline)	491.78		-	

874 of the algorithm sometimes occurred. On the other hand,
875 estimation with the M-quantile approach was smoother but
876 the interpretation of the estimated parameters is more
877 difficult.

878 Finally, Fig. 12 presents the estimated effects obtained
879 from M-quantile and quantile-mixed regression models by
880 quantile for each explanatory variable that is considered in
881 the model. In particular, the solid line represents the pro-
882 posed semiparametric M-quantile random effect regression
883 model and the dashed line stands for an additive quantile

884 regression model (Geraci 2018) which includes a bivariate
885 spline to capture the spatial structure as well as random
886 effects to account for the hierarchy of the data (fitted by the
887 R package *aqmm*). Note that we only plot the point esti-
888 mates (without the point-wise 95% confidence intervals) in
889 order to avoid an overload of Fig. 12. The results confirm
890 that the results based on both models are in the same
891 direction.
892

Fig. 12 Estimated coefficients of quantile regressions (dashed line) and M-quantile regressions (solid line): **a** intercept, **b** fault distance, **c** floor material, **d** wall material, **e** year from construction/last renovation, **f** single building, **g** not in contact with the ground, **h** air conditioning system



Author Proof

References

- Alfó M, Ranalli M, Salvati N (2017) Finite mixtures of quantiles and m-quantile models. *Stat Comput* 27:547–570
- Apte M, Price P, Nero A, Revzan K (1999) Predicting new hampshire indoor radon concentrations from geologic information and other covariates. *Environ Geol* 37:181–194
- AQ3** Bianchi A, Fabrizi E, Salvati N, Tzavidis N (2018) Estimation and testing in M-quantile regression with applications to small area estimation. *Int Stat Rev* 1–30
- Borgoni R (2011) A quantile regression approach to evaluate factors influencing residential indoor radon concentration. *Environ Model Assess* 16:239–250
- Borgoni R, Bianco PD, Salvati N, Schmid T, Tzavidis N (2018) Modelling the distribution of health-related quality of life of advanced melanoma patients in a longitudinal multi-centre clinical trial using m-quantile random effects regression. *Stat Methods Med Res* 27:549–563
- Borgoni R, Quatto P, Soma G, de Bartolo D (2010) A geostatistical approach to define guidelines for radon prone area identification. *Stat Methods Appl* 19:255–276
- Borgoni R, Tritto V, Bigliotto C, de Bartolo D (2011) A geostatistical approach to assess the spatial association between indoor radon concentration, geological features and building characteristics: the Lombardy case, Northern Italy. *Int J Environ Res Public Health* 8:1420–1440
- Bosch RJ, Ye Y, Woodworth GG (1995) A convergent algorithm for quantile regression with smoothing splines. *Comput Stat Data Anal* 19(6):613–630
- Breckling J, Chambers R (1988) M-quantiles. *Biometrika* 75(4):761–771
- Cade B, Noon BR, Flather CH (2005) Quantile regression reveals hidden bias and uncertainty in habitat models. *Ecology* 86:786–800
- Chaudhuri P (1991) Global nonparametric estimation of conditional quantile functions and their derivatives. *J Multivar Anal* 39(2):246–269
- Cinelli G, Tondeur F, Dehandschutter B (2011) Development of an indoor radon risk map of the Walloon region of Belgium, integrating geological information. *Environ Earth Sci* 62:809–819
- Darby S, Hill D, Auvinen A, Barros-Dios J, Baysson J, Bochicchio F, Deo H, Falk R, Forastiere F, Hakama M, Heid I, Kreienbrock L, Kreuzer M, Lagarde F, MSkelSinen I, Muirhead C, Oberaigner W, Pershagen G, Ruano-Ravina A, Ruosteenoja E, Rosario AS, Tirmarche T, Tomsek L, Whitley E, Wichmann H, Doll R (2005) Radon in homes and risk of lung cancer: collaborative analysis of individual data from 13 European case-control studies. *Br Med J* 330(6485):223–226
- Fellner WH (1986) Robust estimation of variance components. *Technometrics* 28(1):51–60
- Fontanella L, Ippoliti L, Sarra A, Valentini P, Palermi S (2015) Hierarchical generalised latent spatial quantile regression models with applications to indoor radon concentration. *Stoch Environ Res Risk Assess* 29:357–367
- Foxall R, Baddeley A (2002) Nonparametric measures of association between a spatial point process and a random set, with geological applications. *J R Stat Soc Ser C* 51(2):165–182
- Gates A, Gundersen L (1992) Geologic controls on radon. Geological Society of America, Washington, DC (Special Paper 271)
- AQ4** Geraci M (2018) Additive quantile regression for clustered data with an application to children's physical activity. ArXiv e-prints
- Geraci M, Bottai M (2014) Linear quantile mixed models. *Stat Comput* 24(3):461–479
- Green B, Miles J, Bradley E, Rees D (2002) Radon atlas of England and Wales. Report nrpb-w26, Chilton NRPB
- Gunby J, Darby S, Miles J, Green B, Cox D (1993) Indoor radon concentrations in the United Kingdom. *Health Phys* 64:2–12
- Huber P (1981) Robust statistics. Wiley, New York
- Huggins RM (1993) A robust approach to the analysis of repeated measures. *Biometrics* 49(3):715–720
- Huggins RM, Loesch DZ (1998) On the analysis of mixed longitudinal growth data. *Biometrics* 54(2):583–595
- Hunter N, Muirhead C, Miles J, Appleton JD (2009) Uncertainties in radon related to house-specific factors and proximity to geological boundaries in England. *Radiat Prot Dosim* 136:17–22
- Jacobi W (1993) The history of the radon problem in mines and homes. *Ann ICRP* 23(2):39–45
- Jones M (1994) Expectiles and m-quantiles are quantiles. *Stat Probab Lett* 20:149–153
- Kaufman L, Rousseeuw P (1990) Finding groups in data: an introduction to cluster analysis. Wiley, New York
- Kemski J, Klingel R, Siehl A, Valdivia-Manchego M (2009) From radon hazard to risk prediction-based on geological maps, soil gas and indoor measurements in Germany. *Environ Geol* 56:1269–1279
- Koenker R (2005) Quantile regression. Cambridge University Press, New York
- Koenker R, Bassett G (1978) Regression quantiles. *Econometrica* 46:33–50
- Koenker R, Mizera I (2004) Penalized triograms: total variation regularization for bivariate smoothing. *J R Stat Soc Ser B* 66(1):145–163
- Koenker R, Ng P, Portnoy S (1994) Quantile smoothing splines. *Biometrika* 81(4):673–680
- Kreienbrock L, Kreuzer M, Gerken M, Dingerkus M, Wellmann J, Keller G, Wichmann H (2001) Case-control study on lung cancer and residential radon in western Germany. *Am J Epidemiol* 89(4):339–348
- Krewski D, Lubin MAJH, Zielinski JM, Catalan V, Field R, Klotz J, Letourneau E, Lynch C, Lyon J, Sandler D, Schoenberg D, Steck J, Stolwijk C, Weinberg C, Wilcox H (2005) Residential radon and risk of lung cancer: a combined analysis of seven North American case-control studies. *Epidemiology* 16(4):137–145
- Levesque B, Gauvin D, McGregor R, Martel R, Gingras S, Dontigny A, Walker W, Lajoie P, Levesque E (1997) Radon in residences: influences of geological and housing characteristics. *Health Phys* 72:907–914
- Lubin J, Boice J (1997) Lung cancer risk from residential radon: a meta-analysis of eight epidemiological studies. *J Natl Cancer Inst* 89(1):49–57
- Nero A, Schwehr M, Nazaroff W, Revzan K (1986) Distribution of airborne radon-222 concentrations in US homes. *Science* 234:992–997
- Newey WK, Powell JL (1987) Asymmetric least squares estimation and testing. *Econometrica* 55(4):819–847
- Opsomer J, Claeskens G, Ranalli M, Kauermann G, Breidt F (2008) Nonparametric small area estimation using penalized spline regression. *J R Stat Soc Ser B* 70(1):265–283
- Organization WH (2009) WHO handbook on indoor radon: a public health perspective. WHO Library Cataloguing-in-Publication Data
- Pratesi M, Ranalli M, Salvati N (2009) Nonparametric m-quantile regression using penalized splines. *J Nonparametr Stat* 21:287–304
- Price P, Nero A, Gelman A (1996) Bayesian prediction of mean indoor radon concentrations for Minnesota counties. *Health Phys* 71:922–936
- R Core Team (2017) R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna

- 1022 Rowe J, Kelly M, Price L (2002) Weather system scale variation in
1023 radon-222 concentration of indoor air. *Sci Total Environ*
1024 284:157–166
- 1025 Ruppert D, Wand M, Carroll R (2003) *Semiparametric regression*.
1026 Cambridge University Press, Cambridge
- 1027 Sarra A, Fontanella L, Ippoliti L, Valentini P, Palermi S (2016)
1028 Quantile regression and Bayesian cluster detection to identify
1029 radon prone areas. *J Environ Radioact* 164:354–364
- 1030 Shi X, Hoftiezer D, Duell E, Onega T (2006) Spatial association
1031 between residential radon concentration and bedrock types in
1032 New Hampshire. *Environ Geol* 51:65–71
- 1033 Smith B, Field R (2007) Effect of housing factor and surficial uranium
1034 on the spatial prediction of residential radon in Iowa. *Environ-*
1035 *metrics* 18:481–497
- 1036 Smith B, Zhang L, Field R (2007) Iowa radon leukemia study: a
1037 hierarchical population risk model. *Stat Med* 10:4619–4642
- 1038 Sundal A, Henriksen H, Soldal O, Strand T (2004) The influence of
1039 geological factors on indoor radon concentrations in Norway. *Sci*
1040 *Total Environ* 328:41–53
- Tiefelsdorf M (2007) Controlling for migration effects in ecological
disease mapping of prostate cancer. *Stoch Environ Res Risk*
Assess 21:615–624
- Tzavidis N, Salvati N, Schmid T, Flouri E, Midouhas E (2016)
Longitudinal analysis of the strengths and difficulties question-
naire scores of the millennium cohort study children in England
using m-quantile random effects regression. *J R Stat Soc Ser A*
179(2):427–452
- USEPA (1992) National residential radon survey: summary report.
Technical Report EPA/402/R-92/011, United States Environ-
mental Protection Agency, Washington, DC
- Wang Y, Lin X, Zhu M, Bai Z (2007) Robust estimation using the
Huber function with a data dependent tuning constant. *J Comput*
Graph Stat 16(2):468–481
- Yu K, Lu Z, Stander J (2003) Quantile regression: applications and
current research areas. *Statistician* 52(3):331–350

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

UNCORRECTED PROOF