

# Environment-Aware Regression for Indoor Localization based on WiFi Fingerprinting

**Germán Mendoza-Silva**

Institute of New Imaging Technologies, Universitat Jaume I, 12071 Castellón, Spain.

**Ana Cristina Costa**

NOVA IMS, Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisboa, Portugal.

**Joaquín Torres-Sospedra**

UBIK Geospatial Solutions, 12006 Castellón, Spain. (e-mail: jtorres@uji.es)

**Marco Painho**

NOVA IMS, Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisboa, Portugal.

**Joaquín Huerta**

Institute of New Imaging Technologies, Universitat Jaume I, 12071 Castellón, Spain.

**This is the accepted version of the article published in *IEEE IEEE Sensors Journal***

**How to cite:** Mendoza-Silva, G., Costa, A. C., Torres-Sospedra, J., Painho, M., & Huerta, J. (2021). Environment-Aware Regression for Indoor Localization based on WiFi Fingerprinting. *IEEE Sensors Journal*, 1-10. <https://doi.org/10.1109/JSEN.2021.3073878>

## **Funding:**

G. M. Mendoza-Silva gratefully acknowledges funding from grant PREDOC/2016/55 by Universitat Jaume I. J. Torres-Sospedra gratefully acknowledges funding from Ministerio de Ciencia, Innovación y Universidades (PTQ2018-009981).

© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or

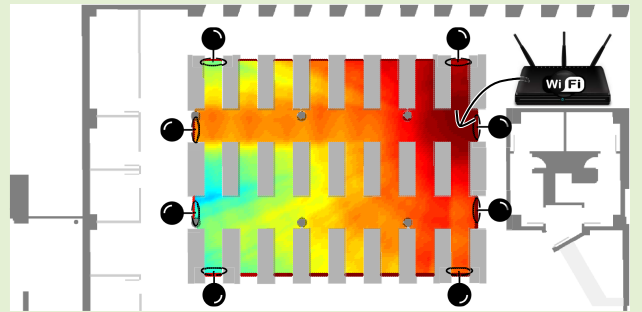
*redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.*

# Environment-Aware Regression for Indoor Localization Based on WiFi Fingerprinting

Germán Martín Mendoza-Silva<sup>ID</sup>, Ana Cristina Costa<sup>ID</sup>, Joaquín Torres-Sospedra<sup>ID</sup>,  
Marco Painho, and Joaquín Huerta<sup>ID</sup>

**Abstract**—Data enrichment through interpolation or regression is a common approach to deal with sample collection for Indoor Localization with WiFi fingerprinting. This paper provides guidelines on where to collect WiFi samples and proposes a new model for received signal strength regression. The new model creates vectors that describe the presence of obstacles between an access point and the collected samples. The vectors, the distance between the access point and the positions of the samples, and the collected, are used to train a Support Vector Regression. The experiments included some relevant analyses and showed that the proposed model improves received signal strength regression in terms of regression residuals and positioning accuracy.

**Index Terms**—Indoor positioning, WiFi fingerprinting, WiFi samples collection, RSS regression.



## I. INTRODUCTION

THE demand for Indoor Positioning Systems (IPS) has already driven academic and commercial research, it is expected that it will dramatically rise in the years to come [1]. Despite the large diversity on related positioning technologies for indoor scenarios, WiFi is one of the most often used. Smartphones and applications relying on Location Based Services (LBS) made WiFi a cost-less approach at the expense of positioning errors around a few meters [2].

Fingerprinting is commonly used with WiFi to provide position indoors. A WiFi fingerprint is a vector with the Received

Signal Strength (RSS) of each WiFi access point (AP) detected in a given position and time. It requires a calibration stage, where samples are collected at well-known positions to create a reference dataset (radio map). In the operational stage, given a new fingerprint measured at an unknown position, the fingerprint method usually provides the centroid of the most similar reference fingerprints as position estimate [3].

Samples collection is known as one of the main challenges of WiFi fingerprinting [4], given that the collection effort can be significant for large areas. The literature suggests to reduce the required effort either by crowdsourcing the collection to volunteers [5], estimating the RSS values applying a propagation model, or applying an interpolation technique to densify an initial reduced radio map [4], [6], [7]. Despite being very valuable, the reliability of position tags and the improper distribution sample position are usual concerns with crowdsourced signal data [8].

This paper addresses the radio map enrichment by applying regression techniques on a proper signal characterization of the environment. Also, through experiments performed on two publicly available databases, we address the problem of choosing the most convenient positions for collecting WiFi fingerprints for radio map creation. Furthermore, we evaluate a new model that applies environment knowledge to Support Vector Regression (SVR), which improves the regression estimates corresponding to extrapolation points in comparison

Manuscript received December 21, 2020; revised March 13, 2021 and March 15, 2021; accepted April 14, 2021. The work of Germán Martín Mendoza-Silva was supported by the Universitat Jaume I under Grant PREDOC/2016/55. The work of Joaquín Torres-Sospedra was supported by the Ministerio de Ciencia, Innovación y Universidades under Grant PTQ2018-009981. The associate editor coordinating the review of this article and approving it for publication was Mr. Francesco Potorti. (Corresponding author: Joaquín Torres-Sospedra.)

Germán Martín Mendoza-Silva and Joaquín Huerta are with the Institute of New Imaging Technologies, Universitat Jaume I, 12071 Castellón de la Plana, Spain (e-mail: gmendoza@uji.es; huerta@uji.es).

Ana Cristina Costa and Marco Painho are with the Nova School of Information Management, Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisbon, Portugal (e-mail: cristina@novaims.unl.pt; painho@novaims.unl.pt).

Joaquín Torres-Sospedra is with UBIK Geospatial Solutions, 12006 Castellón de la Plana, Spain (e-mail: torres@ubikgs.com).

Digital Object Identifier 10.1109/JSEN.2021.3073878

to other extrapolation work shown in WiFi positioning literature.

The main contributions of this paper can be summarized as follows: i) a novel regression model aware of the environment features; ii) a comprehensive analysis of reference position selection to build effective radio maps; and iii) validation in a real-world scenario independent to the research objectives.

## II. BACKGROUND AND RELATED WORKS

A WiFi Access Point (AP) is a networking device that broadcasts one or more wireless networks. A set of RSS values from available APs measured at a specific location throughout a short time interval is called a fingerprint, which can be used for positioning as described in Section I. The quality of the radio map depends on the location of the reference points, the reference point density, the number of samples of each reference points, among many other parameters [9], [10].

However, collecting samples for a radio map requires a notable amount of time [11]. To tackle this problem, two alternatives are usually considered: crowdsourcing and sparse collection. Crowdsourcing has been praised for radio map collection and update [12] at the expense of suffering from low quality of position tags or uneven distributions of the collected samples, whereas sparse collection reduces the collection efforts at the expense of poorer characterizations of the environments. The later approach (sparse collection with regression, interpolation and/or extrapolation models) has been applied to synthetically enrich the radio map for more than 15 years [13], [14], and methods fall in one of the next groups:

- *Sparse recovery* includes, for example, compressed sensing using Singular Value Decomposition (SVD) [15], and radio map interpolation using sparse recovery [6].
- *Interpolation methods* includes traditional interpolation methods [16]–[19]; methods capable of delivering both interpolation and extrapolation like Nearest Neighbor and Inverse Distance Weighting (IWD) [20]; and other interpolation heuristics [21].
- *Extrapolation methods* applied variants based on log-distance path loss model [21]–[23]; on the ray tracing model [24], [25] or the radiosity model [26]–[28].
- *Regression methods* largely includes the application of Gaussian Process Regression (GPR) [29]–[33], although others have also applied Kriging [14], [34]–[36], Geography Weighted Regression (GWR) [37] and Support Vector Regression (SVR) [38].

It is common that radio map enrichment works provide the proportions between points used for fitting and those used for estimations. Talvitie *et al.* [20] concluded that the positions where samples are selected were more important than how many of them were selected. Khalajmehrabadi *et al.* [6] suggested a random selection of reference points and discourage a uniform placement of those points. Ezpeleta *et al.* [16] supported the division in zones arguing that a zone with higher quality of RF signals than other zones required less training points. The importance of the distribution of samples for radio map construction is almost intuitive and acknowledged [39]. However, some works perform random selection

of sample positions for radio map construction [6], [23], [32]. Kanaris *et al.* [40], determined the sample size given a small preliminary set of measurements, suggesting to randomly choose positions from a grid in the number determined by the sample size calculation.

Some radio map enrichment solutions have considered the environment's influence on the signals intensities. The interpolation in Bong and Kim [41] preserved signals discontinuity over the wall. Ali *et al.* [23] used a path loss with wall attenuation factor that introduced an image to count the number of interfering walls. Moghtadaiee *et al.* [21] fitted a log-distance model independently for each architectural zone and created an interpolation that considered only sample at similar distances to the target AP. Some authors [14], [34]–[36] used Kriging, but only considered the Euclidean distance for describing the spatial dependency, which does not hold true for indoor environments. [39] fitted a log distance path loss model for each target position, giving to the samples used for fitting distinct weights (using a kernel density estimation) based on their distances to the target position. Du *et al.* [37] applied GWR, which computed several local models instead a single global one. They used the distance between the emitters and the sample points as predictor variables.

The distribution of samples necessarily should take the layout of the environment into account, not only regarding where it is possible to collect samples, but where is convenient to collect them. The indoor environments strongly influence the WiFi and BLE signals, and the decision on the collection distribution should be aware of it. The radio map enrichment method should ideally be also aware of the target environment, i.e., of the obstacles and the positions of the emitters.

## III. MATERIALS AND METHODS

### A. Selected Datasets

This work is built on top of two public WiFi fingerprinting datasets: the Library dataset [42] and the Mannheim dataset [43]. Partial versions of both datasets will be used to analyse the influence of position distribution and the influence of AP strength on position accuracy. Moreover, they will be used to analyse the influence of AP strength on RSS regressions. For the evaluation of our proposed environment-aware regression model only the Library dataset will be used.

The Library dataset was collected in two floors of the Library building of University Jaume I (UJI) and the systematic data collection was repeated multiple times in a time span of 25 months. There are six WiFi fingerprints per each reference point and each of the two directions at which the collection subject was facing. Also, as the data contained information about a 620 AP, a selection of the 52 most relevant APs was performed (as done in Torres-Sospedra *et al.* [44]) to ease the analyses and reduce the noise created by a large number of intermittent APs. The collection area is a relatively small environment that covered about  $15 \times 10$  m. The average distance between reference points is about 2 m.

The Mannheim dataset was collected in the Mannheim University. The collection area comprises a medium-scale environment, covering about  $50 \times 36$  m of corridors of a university department. The fingerprints are on a 1.5 m grid [43], [45]

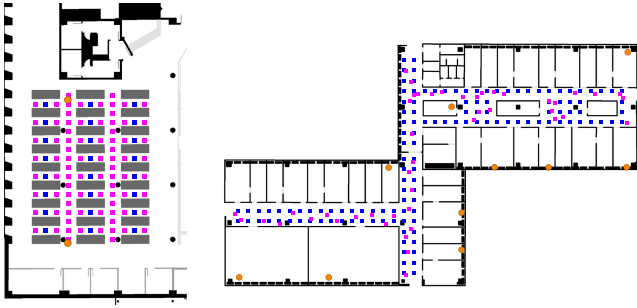


Fig. 1. Library 5<sup>th</sup> floor (left) and mannheim (right) floormaps. Blue and magenta dots represent training and test reference points, respectively. APs positions are drawn with orange circles. Other APs may lay out of the areas.

166 and the positions of 10 APs are known. The dataset contained  
 167 110 fingerprints per reference point. Out of 110 samples,  
 168 we randomly selected 10 to ease the analyzes and have a  
 169 number of samples that is closer to that of the Library dataset.  
 170 Both the original Mannheim and the Library datasets provided  
 171 their position tags using a local coordinate system that allows  
 172 distance computation using the Euclidean distance.

173 Figure 1 shows the operational area of the two evaluation  
 174 environments. The structural barriers were manually created  
 175 from floor plans. Thick walls were drawn in black color and  
 176 thin walls were drawn with a light shade of gray in the  
 177 image, whose intensity values are used by eq.(2). Figure 1 also  
 178 presents the distribution of training and test reference points,  
 179 as well as the position of some APs. The higher the density  
 180 of APs and reference points in the operational area, the lower  
 181 expected positioning error. In both cases, some APs lay out  
 182 the floormap or have an unknown location.

### 183 B. Environment-Aware Regression on WiFi Radio Maps

184 This work presents a regression model that integrates the AP  
 185 reference position and a floor plan of the area. The reference  
 186 position is used as a raw indication of where the AP is. The  
 187 position of APs inside or very close to the collection area  
 188 can be determined with, for instance, the weighted centroid  
 189 or the method proposed in [46]. The approximate position  
 190 of an AP can be also manually obtained by measuring the  
 191 signal intensity with a smartphone application walking in the  
 192 area. However, the accuracy for AP location is low for those  
 193 APs that are away from the operational area and an indicator  
 194 of the relative direction is obtained instead. Those far APs  
 195 are typically detected with a maximum intensity weaker than  
 196  $-60$  dBm. Determining whether an AP is within the collection  
 197 area could be done, for instance, using the Situation Goodness  
 198 test presented in [46] if a relatively dense sample collection is  
 199 available.

200 Figure 2 introduces an example in the Library environment  
 201 (5<sup>th</sup> floor). It shows the mean RSS values per reference point  
 202 for 3 APs, which will later be used to evaluate the proposed  
 203 regression model. The APs with IDs 15 and 49 are inside the  
 204 collection area. Their positions shown in the figure are about  
 205 half a meter and more than a meter away from the actual  
 206 device positions, respectively. The position of the device that

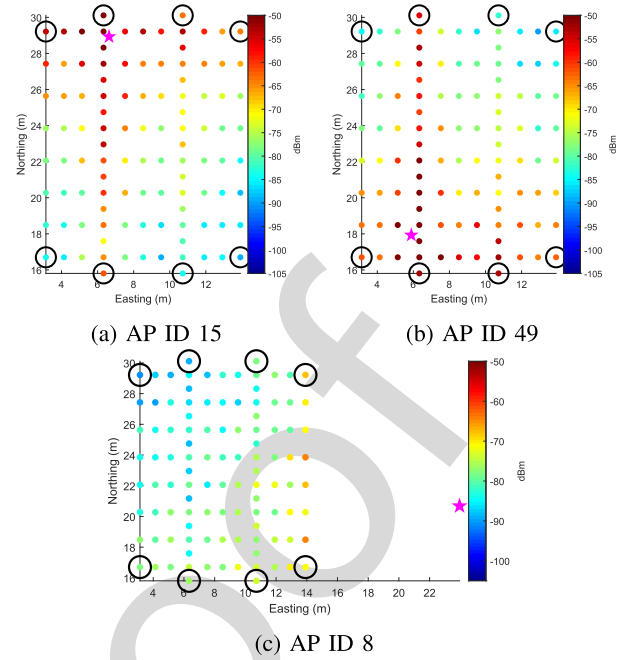


Fig. 2. Mean RSS values per reference point and device reference positions of three APs (Library, 5<sup>th</sup> floor). The device position is indicated with a star. Circles highlight the reference points whose values were used to train regression models.

207 emitted the AP with ID 8 was unknown. The position shown  
 208 in the figure is anyway a useful estimation of the actual AP  
 209 direction.

210 In the proposed model, the predictor variables include  
 211 the target point's position components, the AP's reference  
 212 position and information from the environment floor map.  
 213 Moreover, we applied a data transformation before and after  
 214 the application of the regression method, so that the values  
 215 of the response variable are determined as  $\log_{10}(-RSS)$  (as  
 216 a distance indicator) and the RSS estimate is computed as  
 217  $-(10^{est})$  if  $est$  is an estimate provided by the regression  
 218 model. The positions of points used for training and testing  
 219 the model are expressed in the local coordinate system. Thus,  
 220 their coordinates need to be transformed into image coordi-  
 221 nates (cell positions or *pixels*) before applying the proposed  
 222 model. The following definitions assume positions in image  
 223 coordinates (i.e. *pixels* not meters).

224 Let  $rp = (rp_x; rp_y)$  be the position of a reference  
 225 point used for training the model. Let  $ap = (ap_x; ap_y)$  be  
 226 the position of the AP targeted for regression. Let  $B_{rp} =$   
 227  $\{(x_1, y_1), \dots, (x_k, y_k)\}$  be the line that connects  $rp$  and  
 228  $ap$ . The cell positions that constitute the line are determined  
 229 using the Bresenham's line algorithm [47]. The values of predictor  
 230 variables for  $rp$  are:

$$231 P_{rp} = \{rp_x, rp_y, \frac{d_{rp} + 1}{2}, F_{rp}\}, \quad (1)$$

232 where  $F_{rp} = \{f_1, \dots, f_k, \dots, f_n\}$  and  $d_{rp}$  is the Euclidean  
 233 distance between  $rp$  and  $ap$ . The value  $f_i$  is computed as:

$$234 f_i = \begin{cases} \log_2(2 + 255 - Im(x_i, y_i)) & \text{for } 1 \leq i \leq k \\ 0 & \text{for } k < i \leq n \end{cases} \quad (2)$$



235 where  $x_i$  and  $y_i$  are the position components of the  $i^{th}$  point in  
 236  $B_{rp}$ ,  $Im$  is the image representation of the environment, and  
 237  $Im(x_i, y_i)$  is the cell value in the image  $Im$  whose position  
 238 is  $(x_i, y_i)$ . The value of  $n$  is the maximum number of points  
 239 that may have a line connecting the positions of the AP and a  
 240 point in the environment representation. If  $ap$  lies beyond the  
 241 environment represented by  $Im$ , the image is enlarge applying  
 242 a padding of zeros. In other words,  $Im(x, y) = 0$  for all  $(x, y)$   
 243 that lies beyond the environment representation.

---

**Algorithm 1** : Regression Model Training for an AP
 

---

**Input:**  $Im, RPL, SI, apl$  **Output:** The trained regression model  $M$

- 1 Compute  $ap = (ap_x, ap_y)$ , the position of  $apl_j$  in  $Im$
  - 2 **for** each  $rpl_j$  in  $RPL$  **do**
  - 3     Get  $rp = (rp_x, rp_y)$ , the position of  $rpl_j$  in  $Im$
  - 4     Get  $B_{rp}$ , as stated previously
  - 5     Get  $P_{rp}$ , as stated in Equation 1
  - 6     Set  $p_j = P_{rp}$
  - 7     Get  $resp_j = \log_{10}(-si_j)$
  - 8 **end**
  - 9 Build  $M$  by training SVR using  $\{p_j\}$  as predictors data and  $\{resp_j\}$  as responses data
- 

---

**Algorithm 2** : Signal Prediction for an AP
 

---

**Input:**  $Im, TPL, apl, M$

**Output:** The predicted intensities  $SO$

- 1 Compute  $ap = (ap_x, ap_y)$ , the position of  $apl_j$  in  $Im$
  - 2 **for** each  $tpl_j$  in  $TPL$  **do**
  - 3     Get  $rp = (rp_x, rp_y)$ , the position of  $tpl_j$  in  $Im$
  - 4     Get  $B_{rp}$ , as stated previously
  - 5     Get  $P_{rp}$ , as stated in Equation 1
  - 6     Set  $p_j = P_{rp}$
  - 7     Get  $est_j$  using  $M$  with  $\{p_j\}$  as predictors values
  - 8     Set  $so_j = -(10^{est_j})$
  - 9 **end**
  - 10 Set  $SO = \{so_j\}$
- 

244 Algorithm 1 resumes the process of training the proposed  
 245 regression model. Its inputs are the environment image  $Im$ ,  
 246 the positions (expressed in a local coordinate system) of  
 247 collection points  $RPL = \{rpl_j\}$  and their respective RSS  
 248 values  $SI = \{si_j\}$  measured for an AP. Once the model  $M$  is  
 249 ready, it serves for predicting the RSS values  $SO = \{so_j\}$  for  
 250 a set of positions  $TPL = \{tpl_j\}$  using the Algorithm 2.

251 The set  $F_{rp}$  in Equation 1 is a representation of the obstacles  
 252 between  $rp$  and  $ap$  using the information of the image's cells  
 253 that lie in that path. The cell values in the image  $Im$  represent  
 254 either free space or an obstacle (black or white). Thus,  
 255 the model is trained to learn the influence of an obstacle cell  
 256 value at a given distance from an AP in the signal propagation.  
 257 This work did not differentiate among distinct types of obstacle  
 258 materials for simplicity, despite Equation 2 allows the range  
 259  $[1, \dots, 255]$  for obstacle representation. Setting appropriate  
 260 opaqueness for each material requires additional consideration

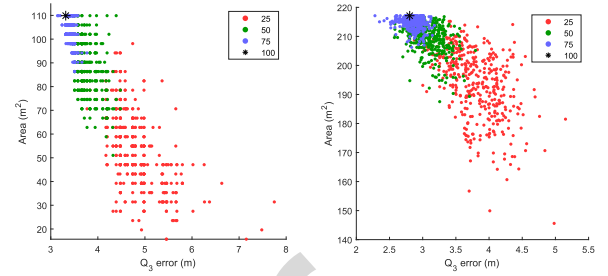


Fig. 3. Relation between training covered area and positioning accuracy for Library 5th Floor (left) and Mannheim (right).

and measurements. Equation 1 includes half of the distance  
 between  $rp$  and  $ap$ . Using the actual value of the distance  
 significantly decreased the obstacles influences in the model.  
 The number of variables presented in Equation 1 depends on  
 the environment and the AP position. Finally, according to our  
 experience, we selected the Support Vector Regression (SVR)  
 with a linear kernel function as regressor.

## IV. EXPERIMENTS AND RESULTS

### A. Influence of RPs Distribution on IPS Accuracy

The goal of the radio map in WiFi fingerprinting is to  
 characterize the signal propagation in the target environment.  
 As the main fingerprint methods (including  $k$ -NN) can only  
 provide position estimates within the convex hull of the  
 reference sample locations, we hipotetise that the number  
 and distribution of the collected samples are strongly related  
 quality of the radio map and, hence, the accuracy of the IPS.

For that purpose, we have evaluated the performance of  
 the radio map in two environments and four different cases:  
 with 100%, 75%, 50% and 25% of RPs. Except for 100%,  
 we repeated the evaluation 400 times with different initializa-  
 tion to cover multiple random scenarios. In all cases, we report  
 the results provided by the optimal  $k$ -value (from the set  
 $[1, \dots, 15]$ ). The results are reported as a scatter in Figure 3  
 for Library 5<sup>th</sup> floor (left) and Mannheim (right).

Every point in the figure represent the area of the reduced  
 radio map's convex hull and the best accuracy reported by the  
 $k$ -NN method with that data set. The accuracy corresponds  
 to the  $Q_3$  value, i.e., the 75<sup>th</sup> percentile as done in IPIN  
 Competition. The color indicates the size of the radio map case  
 (100%, 75%, 50% and 25%). A clear trend can be observed  
 in the two environments, the large the area covered, the best  
 positioning accuracy. In contrast, the worst positioning results  
 came when the convex hull of the reference radio map was  
 small. This is because the kNN method can provide position  
 estimates only within the convex hull of the reference points.  
 Good accuracy can be reached with a reduced radio map if  
 the reference points cover the full operational area.

The figure also shows that the distribution of reference  
 points is relevant. Even for a high covered area, the positioning  
 accuracy can vary up to more than 2 m in the three cases.  
 The largest differences in positioning are observed for cases  
 with low RPs density (i.e. 25%). To evaluate the relation  
 between covered area and accuracy, we calculated the Pearson  
 correlation between the area and the Positioning error in the

third quartile in the 1201 points. The correlation factor ( $\rho$ ) for Mannheim is  $-0.77$ , whereas it is  $-0.89$  for Library. In both cases, the significance ( $p$ -value) is much lower than 0.05 showing that the inverse correlation is statistically significant.

Our hypothesis is that placing reference points near the inner boundary of the collection area would maximize the covered area and assure that test positions are located inside the convex of the training positions. Finding those positions is a trivial task and can be provided by, for instance, alpha-shape [48]. Thus empirical data collection can be optimized to relevant places according to the imposed restrictions. The restrictions will somehow will be an indicator of the density and distribution of the empirical reference points, which will be located only at feasible locations (e.g. there are no samples inside a wall). If the radio map needs to be enriched, regression can be used to synthetically generate new reference samples in those positions that lack of empirically collected data.

One strategy for creating the set of reference points is to first add reference points lying close to environment boundary and later add a number of points  $mp$  that maximize the mean minimum distance among the points in the set. In kNN, the estimated position is commonly computed as the centroid of the positions of the most similar samples in the training dataset. Thus, maximizing the minimum distance among the reference points reduces the areas without position estimates produced by kNN. Such an even distribution of point also benefits regressions as it provides intermediate positions that help explain non-linear behaviors. The value of  $mp$  may be dictated by the affordable collection effort. For low values of  $mp$ , like those below 20, a brute force approach may be applied to determine the  $mp$  positions of the reference points. For large  $mp$  values, a Monte Carlo approach [49] can be used. This work used an optimization approach based on agents moving under repulsion forces [50].

To explore the convenience of using the previous training points distribution, the Pearson correlation test was applied between the mean minimum distance and the positioning error for several distributions of training points. The tests were performed 400 independent times (with random sets of reference points that included the shape boundary) separately for each of the two environments. The position estimations were obtained with kNN, using the best  $k$  for the training set.

Table I presents the correlation results. The negative correlation between the mean minimum distance and the positioning accuracy is not statistically significant. For the Library environment, the negative weak to moderate correlation appears only for large sets, and it is statistically significant for them. The correlation is consistently negative for all set sizes in the Mannheim environment. However, its statistical significance does not show a clear pattern. The results from Table I suggest that the distribution of the inner reference points proposed above is beneficial for environments that are large or have relatively dense collections. Despite it is desirable to avoid the existence of non-positionable zones, alternative distributions may be preferable for other environments.

TABLE I  
CORRELATION ( $\rho$ ) AND STATISTICAL SIGNIFICANCE ( $p$ -VALUE) BETWEEN THE MEAN MINIMUM DISTANCE AMONG TRAINING POINTS AND THE THIRD QUANTILE OF THE POSITIONING ERROR ( $Q_3$ ) FOR DIFFERENT SIZES OF THE RADIO MAP (FROM 25% TO 90% OF RPs)

	Library				Mannheim			
	%RPs	mean k	$\rho$	$p$ -value	$Q_3$ (m)	mean k	$\rho$	$p$ -value
25	2	0.004	0.932	4.740	3	-0.144	0.004	3.898
30	2	0.075	0.135	4.480	3	-0.069	0.167	3.562
35	2	0.122	0.015	4.390	3	-0.120	0.016	3.403
40	2	0.091	0.068	4.171	3	-0.086	0.086	3.319
45	2	0.024	0.628	3.971	3	-0.146	0.003	3.193
50	2	0.060	0.233	3.661	3	-0.131	0.009	3.146
55	2	0.036	0.470	3.576	3	-0.130	0.010	3.071
60	3	-0.036	0.471	3.576	3	-0.077	0.122	2.990
65	3	-0.124	0.013	3.505	3	-0.085	0.088	2.966
70	3	-0.200	0.000	3.466	3	-0.075	0.134	2.926
75	4	-0.313	0.000	3.428	4	-0.152	0.002	2.900
80	4	-0.395	0.000	3.390	4	-0.138	0.006	2.874
85	4	-0.302	0.000	3.322	4	-0.056	0.267	2.864
90	4	-0.279	0.000	3.318	4	-0.131	0.009	2.833

## B. Influence of AP Strength on Positioning Accuracy

It is known that the signal strength from an AP logarithmically decreases as the distance to the AP increases. Thus, it is expected that the closer to the emitter the larger the expected variations in the signals. A radio map should grasp as much as the signal variations in the environment as possible. Having reference points close to the emitter increases the likelihood of incorporating much of those variations.

This subsection explores the correlation between AP proximity to the collection area and the positioning accuracy of a kNN method. Determining the distance to an AP requires knowing the actual position of the AP. Given that the knowledge of AP positions is commonly not assumed for fingerprinting, we inferred proximity from the RSS values. The RSS values for an AP measured in an area should be strong if the AP is close to that area or inside it.

Let us assume a radio map  $RM$  (training set) and a test set. Let  $max_a = \max(\{r_{p,i,a}\})$  be the strongest RSS value for the  $a^{th}$  detected AP in  $RM$ , with  $1 \leq a \leq m$  and  $m$  being the number of APs. Let  $qap$  (median inferred proximity) be the  $Q_2$  value of  $\{max_a\}$ . Let  $qpe$  (positioning accuracy) be the  $Q_3$  value of positioning errors obtained by a kNN method using the above training and test sets.

Here, we also created 400 random subsets containing the 25% of an original training set (either for the Library or Mannheim). For each subset  $RM_s$ , the  $qap_s$  and  $qpe_s$  were computed. For  $qpe_s$ , the kNN method used  $RM_s$  as training set and the original test set. Then, the Pearson correlation test was applied on the sets  $\{qap_s\}$  and  $\{qpe_s\}$ , with  $1 \leq s \leq 400$ . The test results are shown in Table II. For the two environments, the correlation results were statistically significant. The low to moderate negative correlation indicates that high accuracy is associated with low proximity values (weak RSS). Thus, the results suggest the convenience of distributing some reference points in zones of the collection area where nearby APs are may result in large signal variations.

TABLE II

CORRELATION TEST RESULTS BETWEEN  $qap$  (MEDIAN INFERRED PROXIMITY) AND  $qpe$  (POSITIONING ACCURACY)

Environment	$\rho$	$p$ -value
Library	-0.37	$\approx 0$
Mannheim	-0.28	$\approx 0$

TABLE III

CORRELATION BETWEEN MEAN VALUES OF SIGNAL STRENGTH IN THE ENVIRONMENT AND MEAN VALUES OF REGRESSION RESIDUALS

Environment	$\rho$	$p$ -value
Library	0.88	$\approx 0$
Mannheim	0.24	$\approx 0$

### C. Influence of AP Strength on RSS Regressions

399

400 The following experiments addressed the notion of the  
 401 convenience of having more reference points close to nearby  
 402 APs in relation to the regression or interpolation results.  
 403 The goodness of a regression or an interpolation applied to  
 404 radio map densification is normally assessed by the difference  
 405 between the estimated RSS and their actual values. The  
 406 interpolation methods used in the experiments were Nat-  
 407 ural Neighbours [51], (Bi)Cubic Interpolation [52], [53] and  
 408 Inverse Distance Weighting [54]. The regression methods used  
 409 in the experiments were Support Vector Machines (SVM) [55],  
 410 Gaussian Process [56], Generalized Linear Models [57], Deci-  
 411 sion Trees (DT) [58], and Ensembles of Decision Trees [59].  
 412 The interpolation and regression methods, hereinafter only  
 413 called regression methods, were applied using training points  
 414 to fit the model and tests point to compute RSS estimates. The  
 415 mean RSS value for an AP and a reference point was used to  
 416 train the regression model for an AP and to later compute the  
 417 regression residuals. The residuals are the AP-wise absolute  
 418 difference between RSS estimates provided by the regression  
 419 and the actual RSS used for training.

420 Table III shows the correlation values between signal  
 421 strength and regression residuals for each environment. Let  
 422  $S_j = \{s_1, \dots, s_n\}$  and  $R_j = \{r_{j,1}, \dots, r_{j,n}\}$  be two sets, where  
 423  $n$  is the number of APs detected in that environment. The value  
 424  $s_i$  was computed as the mean RSS value of the  $i^{th}$  AP in the  
 425 environment, considering all reference points. The value  $r_{j,i}$   
 426 was computed as the mean of the residual values obtained  
 427 for the  $i^{th}$  AP applying the  $j^{th}$  regression method in the  
 428 environment. The values for the signal strength and regression  
 429 residuals used for the correlation test in an environment are the  
 430 sets  $\{S_1, \dots, S_m\}$  and  $\{R_1, \dots, R_m\}$ , where  $m$  is the number  
 431 of regression methods.

432 The correlation is statistically significant for the two  
 433 environments. The correlation magnitude is weak for the  
 434 Mannheim environment but notable for the Library environ-  
 435 ment. The higher the median value of the signal strength in  
 436 the environment, the larger the residuals of the regressions.  
 437 The correlation difference between the two environments is  
 438 a likely result of the dimensions of the environments. The  
 439 Mannheim environment is large, and thus the detected signal

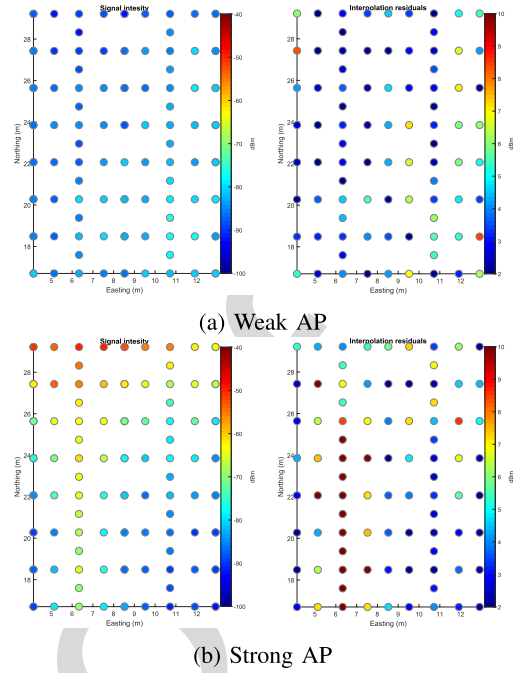


Fig. 4. Mean value of residuals distribution compared to mean value of RSS for the library environment.

440 intensities for an AP can be very strong in some areas and very  
 441 weak at some other areas. Very strong and very weak signal  
 442 intensities are not detected for the same AP in the Library  
 443 environment.

444 Figure 4 shows the relation between the strength with which  
 445 an AP is seen in an environment and the regression goodness.  
 446 The investigation was performed for two APs in the Library  
 447 environment (one with weak and one with strong RSS values).  
 448 The charts from Figure 4 present for each AP includes the  
 449 median value for the RSS values of the AP at each reference  
 450 point and the median value of the regression residuals at  
 451 each reference point. In particular, figure 4a shows regression  
 452 residuals of moderate values for the weak AP, while Figure 4b  
 453 shows regression residuals for the strong AP that are not only  
 454 notably larger than those for the weak AP but also mainly  
 455 situated in a specific zone of the environment.

456 The charts suggest that for weak, far away APs, the regres-  
 457 sion requires only a few samples to train a model, as the  
 458 APs signals are only weakly affected by the environment.  
 459 However, the strength values of signals from APs near the  
 460 target environment heavily depend on the Line of Sight (LOS)  
 461 and Non Line of Sight (NLOS) situations.

462 Table IV presents the spatial auto-correlation test as  
 463 obtained by Moran's I [60] for the two antennas addressed  
 464 in Figure 4. The table suggests that for APs strongly seen  
 465 across the environment the distribution of regression residuals  
 466 is not random and tends to organize in clusters; while for APs  
 467 weakly seen in the environment the distribution of residuals is  
 468 likely random. As stated in the literature [61], the environment  
 469 influence is less significant for weak than for strong signals.  
 470 Furthermore, the signal in free space follows a logarithmic  
 471 decay, i.e., the farther from the AP the slower the decay  
 472 rate. The tested regression models fail to account for a spatial



TABLE IV  
SPATIAL AUTO-CORRELATION (MORAN'S I)  
OF REGRESSION RESIDUALS

Behavior	$Q_2$ of RSS	$Q_2$ of residuals	$z$ -score	$p$ -value
Weak AP	-83	4	0.920	0.357
Strong AP	-74	6	7.702	$\approx 0$

473 process induced by the environment for strong signals. Thus,  
474 samples are required in zones of LOS and NLOS with respect  
475 to nearby APs, given that the RSS values in those two  
476 situations can be significantly different.

477 Given the moderate correlation obtained in some of the  
478 analyses, and that the experiments were only performed in two  
479 environments, a reference point position determination method  
480 is not proposed. However, such determination method may  
481 have the following steps:

- 482 1) Place some reference points in the boundaries.
- 483 2) Distribute the rest of point maximizing the mean mini-  
484 mum distance among reference points.
- 485 3) Adjust the distribution to have some points closer to  
486 nearby APs.
- 487 4) Tend to LOS situations, assuring to place points in LOS  
488 and NLOS situations.

489 This work recommends the previous method steps as a set of  
490 guidelines that follow after the results of the analyses provided  
491 in this section. The most common approach of placing the  
492 reference points on a grid does not take into account the  
493 environment characteristics. The guidelines suggest adapting  
494 the sampling positions to the environment and highlight the  
495 importance of knowing the position of nearby antennas. Thus,  
496 the following two experiments address the environment aware  
497 regression and its evaluation on the Library environment.  
498 We selected the Library as the evaluation environment because  
499 the benefits from including environment knowledge into a  
500 regression model were expected to be greater for the Library  
501 than for Mannheim, as suggested by the correlations shown  
502 in Table III. Furthermore, the Library environment represents  
503 a medium-size open area with many obstacles (bookshelves),  
504 in which a positioning service is commonly desired.

#### 505 D. Environment Aware Regression Assessment

506 The regression models were generated using the reference  
507 points that defined the boundary of the collection area (see  
508 Figure 2), which represent less than 8% of all available  
509 reference points. The remaining reference points were used  
510 to compute the regression residuals. The experiments only  
511 included APs that had measurements for all reference points.

512 Figure 5 presents the regression estimates for APs 15, 49  
513 and 8 provided by a baseline that combines Natural Neighbour  
514 interpolation and Gradient Extrapolation and by the proposed  
515 regression model based on Support Vector Machine. For our  
516 proposed Model, the images were smoothed using 9 pixels  
517 square windows convolution.

518 Given the small number of training points, the two regres-  
519 sions performed remarkably well for the APs located inside  
520 the collection area, APs 15 and 49. The proposed regression

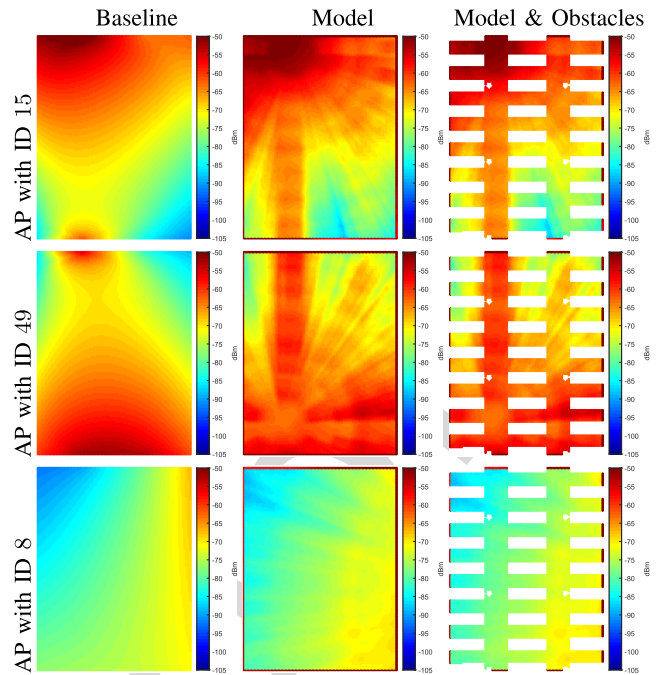


Fig. 5. Regression estimates for APs 15, 49 and 8.

TABLE V  
75<sup>th</sup> PERCENTILE OF REGRESSION RESIDUALS IN DB

AP ID	1	6	8	15	17	49	51	52	54	69
Model	6.7	6.8	4.7	8.8	8.0	7.8	9.5	8.4	8.5	11.5
Baseline	10.7	10.3	6.6	13.6	13.2	12.2	13.2	10.1	11.2	16.2
Difference	4	3.5	1.9	4.8	5.2	4.4	3.7	1.7	2.7	4.7

521 can clearly capture the influence of obstacles in the radio  
522 map. For an AP outside the collection zone, the difference  
523 between Baseline and Model regression is not significant as  
524 the environment has little impact on the propagation of weak  
525 signals. The proposed model captures such behavior, and thus  
526 its estimates mostly depend on the distance to APs.

527 Figure 6 presents the regression residuals obtained using the  
528 baseline and our proposed model. The residuals obtained for  
529 the proposed model are consistently better than those from the  
530 baseline. For AP 15, the maximum residual value was about  
531 10 dBm smaller in the proposed model than in the baseline.  
532 For AP 49, the maximum residual values were similar for  
533 the two approaches. However, the proposed model performed  
534 notably better than the baseline regarding percentiles between  
535 the 25<sup>th</sup> and 75<sup>th</sup>. For AP 8, the difference in residual values  
536 is less notable than for the previous two AP, which is in part  
537 a result of notably lower residual values.

538 Table V presents the 75<sup>th</sup> percentile of regression residuals  
539 for the proposed model and the baseline method. The results  
540 are provided for some relevant APs, i.e., those APs with  
541 valid measurements available for all (106) reference points.  
542 Additionally, we included AP 71 (which had measurements  
543 for 105 points) and one weakly seen AP (AP 8). The  
544 proposed method performs better than the baseline for all  
545 selected APs.

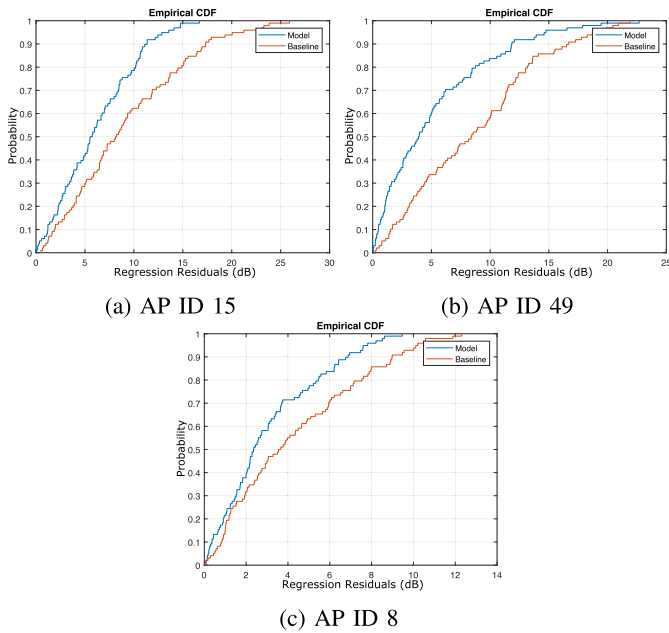


Fig. 6. Regression residuals CDF of baseline and our model.

## V. EMPIRICAL VALIDATION

This section includes the empirical validation by applying together the two main contributions of this paper: the convenient positions where to collect the reference samples and the improved RSS regressor to enhance the radio map. For that purpose, we have used the data collected for the first month from the Library dataset [42]. It corresponds to a real environment with several obstacles (bookshelves and people) whose data collection was independent to this research work.

**Traditional:** The set 01 from the training set was used as reference data (radio map), and the sets 02–05 from the evaluation set were used for evaluation.

**Measurement:** Only the testing data (sets 01–05) was used for reference and evaluation. The 8 points highlighted in Figure 2 are used as training data (radio map), whereas the remaining points are used for evaluation.

**Interpolation – Baseline:** Similar to *Measurement*, but Natural Neighbors interpolation model is applied to increase the density of data in the training set.

**Interpolation – Proposed model:** Similar to *Measurement*, but our proposed interpolation model is applied to increase the density of data in the training set.

Following the ISO 18305 Standard for test and evaluation of localization and tracking systems, we report the results using the mean, median and 95<sup>th</sup> percentile (P95) of the positioning error in Table VI. Additionally, we provide the Third quartile (Q3) as done in the IPIN Competition [62] and the 90<sup>th</sup> percentile (P90).

As expected, the *traditional* approach, where multiple reference positions (24 in this case) are equally distributed in the operational area, is providing the best overall results, except, surprisingly, for the P95 metric. The *measurement* approach (with 8 reference points) is, as expected, providing the worst results as a few reference points are located in the periphery. Both interpolations, the *Natural Neighbors* and

TABLE VI  
RESULTS OF THE EMPIRICAL EVALUATION

Base model.	Mean	Median	Q3	P90	P95
Traditional	3.41	2.83	4.74	6.63	7.94
Measurement	4.26	3.71	5.67	7.98	8.43
Interpolation – Baseline	4.06	3.69	5.67	7.32	8.7
Interpolation – Proposed model	3.94	3.8	5.38	6.82	7.21

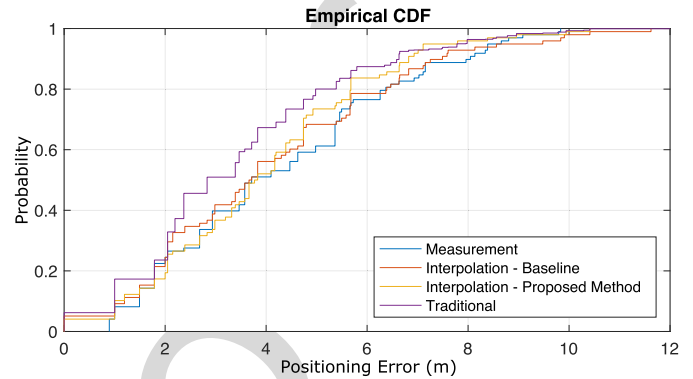


Fig. 7. Positioning accuracy.

our proposed model, improve the results of the *measurement* approach. In general, our model is providing the best results using the reduced set of reference points. With a few reference points, we achieved a mean accuracy below 4 m and percentile errors close to the traditional approach.

Analysing the CDF plot (Figure 7) we can observe that: i) below 30<sup>th</sup> percentile, the traditional approach and both interpolations perform similarly; ii) between 30<sup>th</sup> and 80<sup>th</sup> percentiles, the traditional approach is clearly the best method (at the expense of collecting 3 times more reference data); and iii) the traditional approach and our proposed method have a similar performance in values above 80<sup>th</sup> percentile.

## VI. CONCLUSION

This paper has addressed the reduction of collection efforts for WiFi fingerprinting with two proposals. The first proposal is a set of guidelines to determine convenient positions where to collect WiFi samples. The second proposal is a model that improves the RSS regression estimates for APs that are strongly seen in the collection area. The guidelines were drawn from experiments that analyzed the effect that the distribution of collection points and the intensity of the APs in the environment have in (1) the accuracy of an IPS and (2) in the quality of a regression that could be applied to enrich the radio map. The guidelines highlight the importance of situating collection points around the boundaries of the target environment. Also, zones that are close to APs require more collection points than others. Thus, the position of an AP was shown to be an important piece of information for the determination of collection positions. Furthermore, the regressions and interpolation methods are shown to provide very good estimates for AP weakly seen in the environment.

The proposed model considers the influence of obstacles to improve WiFi RSS regressions for APs strongly seen in the

environment. The model requires an approximate reference position of the AP whose RSS are to be estimated. The reference AP position and raw map information of the obstacles in the environment are used to create the training features for a Support Vector Machine regression. The regression proposal provided RSS estimates better than other regression or interpolation methods in the test environment and selected (strong) APs. The benefits of the regression proposal were also tested according to the positioning accuracy of a kNN method. The kNN was applied (1) using the radio map composed only by collected samples, (2) using the radio map created with other regression or interpolation methods, and (3) using the radio map created with our regression proposal. The best positioning accuracy was obtained using the third option.

The regression model presented in this paper could be considered a first step towards the definition of more general regression models or methods where, for instance, the type of material could be considered. To the best of this work's knowledge, there is no interpolation method, regression method, or tool that allows the direct modeling of the environment influence (presence of obstacles and walls) on a measured phenomenon.

The idea behind the regression model proposed in this paper could inspire others to include the environment characteristics into the existent methods that consider the spatial relation between measurements. We acknowledge that more ambitious conclusions would have reached with a more comprehensive evaluation. However, some methods proposed in the literature are not fully reproducible (some parameters are still missing) and the set of diverse data sets available for positioning do not contain enough information to integrate maps. We, the indoor positioning community, need to adopt and promote reproducible practices as well as creating rich data sets following international standards and ensuring interoperability. Further research is still needed to test the proposed method in a more challenging industrial environments and/or using BLE as positioning technology.

## REFERENCES

- [1] G. M. Mendoza-Silva, J. Torres-Sospedra, and J. Huerta, "A meta-review of indoor positioning systems," *Sensors*, vol. 19, no. 20, p. 4507, Oct. 2019. [Online]. Available: <https://www.mdpi.com/1424-8220/19/20/4507>
- [2] R. Mautz, "Indoor positioning technologies," Ph.D. dissertation, ETH Zürich, Zürich, Switzerland, 2012.
- [3] P. Bahl and V. N. Padmanabhan, "RADAR: An in-building RF-based user location and tracking system," in *Proc. IEEE INFOCOM Conf. Comput. Commun., 19th Annu. Joint Conf. IEEE Comput. Commun. Societies*, Mar. 2000, pp. 775–784, doi: [10.1109/infcom.2000.832252](https://doi.org/10.1109/infcom.2000.832252).
- [4] A. Perez-Navarro *et al.*, "Challenges of fingerprinting in indoor positioning and navigation," in *Geographical and Fingerprinting Data to Create Systems for Indoor Positioning and Indoor/Outdoor Navigation*. Amsterdam, The Netherlands: Elsevier, 2019, pp. 1–20, doi: [10.1016/b978-0-12-813189-3.00001-0](https://doi.org/10.1016/b978-0-12-813189-3.00001-0).
- [5] E. Lohan, J. Torres-Sospedra, H. Leppäkoski, P. Richter, Z. Peng, and J. Huerta, "Wi-Fi crowdsourced fingerprinting dataset for indoor positioning," *Data*, vol. 2, no. 4, p. 32, Oct. 2017. [Online]. Available: <https://www.mdpi.com/2306-5729/2/4/32>
- [6] A. Khalajmehrabadi, N. Gatsis, and D. Akopian, "Structured group sparsity: A novel indoor WLAN localization, outlier detection, and radio map interpolation scheme," *IEEE Trans. Veh. Technol.*, vol. 66, no. 7, pp. 6498–6510, Jul. 2017.
- [7] A. Khalajmehrabadi, N. Gatsis, and D. Akopian, "Modern WLAN fingerprinting indoor positioning methods and deployment challenges," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 3, pp. 1974–2002, 3rd Quart., 2017.
- [8] X. Zhou, T. Chen, D. Guo, X. Teng, and B. Yuan, "From one to crowd: A survey on crowdsourcing-based wireless indoor localization," *Frontiers Comput. Sci.*, vol. 12, no. 3, pp. 423–450, Jun. 2018, doi: [10.1007/s11704-017-6520-z](https://doi.org/10.1007/s11704-017-6520-z).
- [9] D. Zou, W. Meng, S. Han, Z. Gong, and B. Yu, "User aided self-growing approach on radio map construction for wlan based localization," in *Proc. 26th Int. Tech. Meeting Satell. Division Inst. Navigat. (ION GNSS+)*, 2013, pp. 991–997.
- [10] S. Tsruya, R. Dalla Torre, D. Aljadeff, and L. Amir, "Devices, methods, and systems for radio map generation," U.S. Patent 8938255, Jan. 20, 2015.
- [11] F. Meneses, A. Moreira, A. Costa, and M. J. Nicolau, "4—Radio maps for fingerprinting in indoor positioning," in *Geographical and Fingerprinting Data to Create Systems for Indoor Positioning and Indoor/Outdoor Navigation*. New York, NY, USA: Academic, 2019, pp. 69–95. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/B9780128131893000046>
- [12] Y. Zhuang, Z. Syed, J. Georgy, and N. El-Sheimy, "Autonomous smartphone-based WiFi positioning system by using access points localization and crowdsourcing," *Pervas. Mobile Comput.*, vol. 18, pp. 118–136, Apr. 2015. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1574119215000358>
- [13] J. Krumm and J. Platt, *Minimizing Calibration Effort for an Indoor 802.11 Device Location Measurement System*. Redmond, WA, USA: Microsoft Research, Nov. 2003.
- [14] B. Li, Y. Wang, H. K. Lee, A. Dempster, and C. Rizos, "Method for yielding a database of location fingerprints in WLAN," *IEE Proc. Commun.*, vol. 152, no. 5, pp. 580–586, Oct. 2005.
- [15] Z. Gu, Z. Chen, Y. Zhang, Y. Zhu, M. Lu, and A. Chen, "Reducing fingerprint collection for indoor localization," *Comput. Commun.*, vol. 83, pp. 56–63, Jun. 2016. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0140366415003643>
- [16] S. Ezpeleta, J. Claver, J. Pérez-Solano, and J. Martí, "RF-based location using interpolation functions to reduce fingerprint mapping," *Sensors*, vol. 15, no. 10, pp. 27322–27340, Oct. 2015. [Online]. Available: <http://www.mdpi.com/1424-8220/15/10/27322/htm>
- [17] J. Racko, J. Machaj, and P. Brida, "Wi-Fi fingerprint radio map creation by using interpolation," *Procedia Eng.*, vol. 192, pp. 753–758, Jan. 2017. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1877705817326760>
- [18] M. Zhang and W. Cai, "Multivariate polynomial interpolation based indoor fingerprinting localization using Bluetooth," *IEEE Sensors Lett.*, vol. 2, no. 4, pp. 1–4, Dec. 2018.
- [19] L. Xie, X. Jin, M. Zhou, Y. Wang, and Z. Tian, "Cost-efficient BLE fingerprint database construction approach via multi-quadric RBF interpolation," *EURASIP J. Wireless Commun. Netw.*, vol. 2019, no. 1, pp. 1–15, Dec. 2019.
- [20] J. Talvitie, M. Renfors, and E. S. Lohan, "Distance-based interpolation and extrapolation methods for RSS-based localization with indoor wireless signals," *IEEE Trans. Veh. Technol.*, vol. 64, no. 4, pp. 1340–1353, Apr. 2015.
- [21] V. Moghtadaiee, S. A. Ghorashi, and M. Ghavami, "New reconstructed database for cost reduction in indoor fingerprinting localization," *IEEE Access*, vol. 7, pp. 104462–104477, 2019.
- [22] J. S. Seybold, *Indoor Propagation Modeling*. Hoboken, NJ, USA: Wiley, 2005, ch. 9, pp. 208–216.
- [23] M. Ali, S. Hur, and Y. Park, "LOCAL: Calibration-free systematic localization approach for indoor positioning," *Sensors*, vol. 17, no. 6, p. 1213, May 2017. [Online]. Available: <https://www.mdpi.com/1424-8220/17/6/1213>
- [24] A. S. Glassner, *An Introduction to Ray Tracing*. Amsterdam, The Netherlands: Elsevier, 1989.
- [25] M. Ayadi, N. Torjemen, and S. Tabbane, "Two-dimensional deterministic propagation models approach and comparison with calibrated empirical models," *IEEE Trans. Wireless Commun.*, vol. 14, no. 10, pp. 5714–5722, Oct. 2015.
- [26] P. S. Heckhert, "Radiosity in flatland," *Comput. Graph. Forum*, vol. 11, no. 3, pp. 181–192, May 1992.
- [27] M. Cohen, D. Greenberg, D. Immel, and P. Brock, "An efficient radiosity approach for realistic image synthesis," *IEEE Comput. Graph. Appl.*, vol. 6, no. 3, pp. 26–35, Mar. 1986.



- [28] Ó. Belmonte-Fernández, R. Montoliu, J. Torres-Sospedra, E. Sansano-Sansano, and D. Chia-Aguilar, "A radiosity-based method to avoid calibration for indoor positioning systems," *Expert Syst. Appl.*, vol. 105, pp. 89–101, Sep. 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0957417418302112>
- [29] M. M. Atia, A. Noureldin, and M. J. Korenberg, "Dynamic online-calibrated radio maps for indoor positioning in wireless local area networks," *IEEE Trans. Mobile Comput.*, vol. 12, no. 9, pp. 1774–1787, Sep. 2013.
- [30] P. Richter and M. Toledano-Ayala, "Revisiting Gaussian process regression modeling for localization in wireless sensor networks," *Sensors*, vol. 15, no. 9, pp. 22587–22615, Sep. 2015. [Online]. Available: <https://www.mdpi.com/1424-8220/15/9/22587>
- [31] H. Zou, M. Jin, H. Jiang, L. Xie, and C. J. Spanos, "WinIPS: WiFi-based non-intrusive indoor positioning system with online radio map construction and adaptation," *IEEE Trans. Wireless Commun.*, vol. 16, no. 12, pp. 8118–8130, Dec. 2017.
- [32] W. Sun, M. Xue, H. Yu, H. Tang, and A. Lin, "Augmentation of fingerprints for indoor WiFi localization based on Gaussian process regression," *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 10896–10905, Nov. 2018.
- [33] H. Ai, K. Tang, W. Huang, S. Zhang, and T. Li, "Fast fingerprints construction via GPR of high spatial-temporal resolution with sparse RSS sampling in indoor localization," *Computing*, vol. 102, no. 3, pp. 781–794, Mar. 2020, doi: [10.1007/s00607-019-00724-5](https://doi.org/10.1007/s00607-019-00724-5).
- [34] C. Liu, A. Kiring, N. Salman, L. Mihaylova, and I. Esnaola, "A Kriging algorithm for location fingerprinting based on received signal strength," in *Proc. Sensor Data Fusion, Trends, Solutions, Appl. (SDF)*, Oct. 2015, pp. 1–6.
- [35] S.-S. Jan, S.-J. Yeh, and Y.-W. Liu, "Received signal strength database interpolation by Kriging for a Wi-Fi indoor positioning system," *Sensors*, vol. 15, no. 9, pp. 21377–21393, Aug. 2015.
- [36] S. Kram, C. Nickel, J. Seitz, L. Patino-Studencka, and J. Thielecke, "Spatial interpolation of Wi-Fi RSS fingerprints using model-based universal kriging," in *Proc. Sensor Data Fusion, Trends, Solutions, Appl. (SDF)*, Oct. 2017, pp. 1–6.
- [37] Y. Du, D. Yang, and C. Xiu, "A novel method for constructing a WiFi positioning system with efficient manpower," *Sensors*, vol. 15, no. 4, pp. 8358–8381, Apr. 2015. [Online]. Available: <https://www.mdpi.com/1424-8220/15/4/8358>
- [38] N. Hernández, M. Ocaña, J. Alonso, and E. Kim, "Continuous space estimation: Increasing WiFi-based indoor localization resolution without increasing the site-survey effort," *Sensors*, vol. 17, no. 12, p. 147, Jan. 2017. [Online]. Available: <https://www.mdpi.com/1424-8220/17/1/147>
- [39] L. Li, G. Shen, C. Zhao, T. Moscibroda, J.-H. Lin, and F. Zhao, "Experiencing and handling the diversity in data density and environmental locality in an indoor positioning service," in *Proc. 20th Annu. Int. Conf. Mobile Comput. Netw.*, Sep. 2014, pp. 459–470, doi: [10.1145/2639108.2639118](https://doi.org/10.1145/2639108.2639118).
- [40] L. Kanaris, A. Kokkinis, G. Fortino, A. Liotta, and S. Stavrou, "Sample size determination algorithm for fingerprint-based indoor localization systems," *Comput. Netw.*, vol. 101, pp. 169–177, Jun. 2016.
- [41] W. Bong and Y. C. Kim, "Fingerprint Wi-Fi radio map interpolated by discontinuity preserving smoothing," in *Proc. Int. Conf. Hybrid Inf. Technol.* Springer, 2012, pp. 138–145.
- [42] G. Mendoza-Silva, P. Richter, J. Torres-Sospedra, E. Lohan, and J. Huerta, "Long-term WiFi fingerprinting dataset for research on robust indoor positioning," *Data*, vol. 3, no. 1, p. 3, Jan. 2018.
- [43] T. King, S. Kopf, T. Haenselmann, C. Lubberger, and W. Effelsberg, (Apr. 2008). *CRAWDAD Dataset Mannheim/Compass (V. 2008-04-11)*. [Online]. Available: <https://crawdad.org/mannheim/compass/20080411>
- [44] J. Torres-Sospedra, P. Richter, G. Mendoza-Silva, E. S. Lohan, and J. Huerta, "Characterising the alteration in the AP distribution with the RSS distance and the position estimates," in *Proc. Int. Conf. Indoor Positioning Indoor Navigat. (IPIN)*, Sep. 2018, pp. 1–8.
- [45] T. King, T. Haenselmann, and W. Effelsberg, "On-demand fingerprint selection for 802.11-based positioning systems," in *Proc. Int. Symp. World Wireless, Mobile Multimedia Netw.*, Jun. 2008, pp. 1–8.
- [46] G. M. Mendoza-Silva, J. Torres-Sospedra, J. Huerta, R. Montoliu, F. Benitez, and O. Belmonte, "Situation goodness method for weighted centroid-based Wi-Fi aps localization," in *Progress in Location-Based Services*. Springer, 2017, pp. 27–47.
- [47] J. E. Bresenham, "Algorithm for computer control of a digital plotter," *IBM Syst. J.*, vol. 4, no. 1, pp. 25–30, 1965.
- [48] H. Edelsbrunner, D. Kirkpatrick, and R. Seidel, "On the shape of a set of points in the plane," *IEEE Trans. Inf. Theory*, vol. IT-29, no. 4, pp. 551–559, Jul. 1983.
- [49] N. Metropolis, "The beginning of the Monte Carlo method," *Los Alamos Sci.*, vol. 15, no. 584, pp. 125–130, 1987.
- [50] A. Crooks, "Agent-based modeling and geographical information systems," in *Geocomputation, A Practical Primer*. Los Angeles, CA, USA: Sage, 2015, pp. 63–77.
- [51] R. Sibson, "A brief description of natural neighbour interpolation," *Interpreting Multivariate Data*, to be published.
- [52] T. Yang, *Finite Element Structural Analysis*, vol. 2. Englewood Cliffs, NJ, USA: Prentice-Hall, 1986.
- [53] D. Watson, *Contouring: A Guide to the Analysis and Display of Spatial Data*, vol. 10. Amsterdam, The Netherlands: Elsevier, 1992.
- [54] D. Shepard, "A two-dimensional interpolation function for irregularly-spaced data," in *Proc. 23rd ACM Nat. Conf.*, 1968, pp. 517–524.
- [55] V. N. Vapnik, *The Nature of Statistical Learning Theory*. Berlin, Germany: Springer, 1995.
- [56] C. K. Williams and C. E. Rasmussen, *Gaussian Processes for Machine Learning*, vol. 2. Cambridge, MA, USA: MIT Press, 2006.
- [57] J. A. Nelder and R. W. Wedderburn, "Generalized linear models," *J. Roy. Stat. Soc., Ser. A (Gen.)*, vol. 135, no. 3, pp. 370–384, 1972.
- [58] L. Breiman, J. Friedman, C. Stone, and R. Olshen, "Classification and regression trees," in *The Wadsworth and Brooks-Cole Statistics-Probability Series*. Abingdon, U.K.: Taylor & Francis, 1984. [Online]. Available: <https://books.google.es/books?id=JwQx-WOmSyQC>
- [59] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Ann. Statist.*, vol. 29, no. 5, pp. 1189–1232, Oct. 2001.
- [60] P. A. P. Moran, "Notes on continuous stochastic phenomena," *Biometrika*, vol. 37, nos. 1–2, pp. 17–23, 1950.
- [61] A. Goldsmith, "Path loss and shadowing," in *Wireless Communications*. Cambridge, U.K.: Cambridge Univ. Press, 2005, ch. 2, pp. 25–62. [Online]. Available: <https://books.google.es/books?id=n-3ZZ9i0s-cC>
- [62] F. Potorti, A. Crivello, and F. Palumbo, "11—The EvAAL evaluation framework and the IPIN competitions," in *Geographical and Fingerprinting Data to Create Systems for Indoor Positioning and Indoor/Outdoor Navigation (Intelligent Data-Centric Systems)*, J. Conesa, A. Perez-Navarro, J. Torres-Sospedra, and R. Montoliu, Eds. New York, NY, USA: Academic, 2019, pp. 209–224. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/B9780128131893000113>



**Germán Martín Mendoza-Silva** received the bachelor's degree in computer science from the University of Oriente, Cuba, in 2005, the M.Sc. degree in geospatial technologies from WWU, Germany, UNL, Portugal, and UJI, Spain, in 2015, and the Ph.D. degree in informatics from UJI, in 2020. His research interests include WLAN-based indoor positioning, indoor navigation, machine learning, and GIS applications.



**Ana Cristina Costa** received the Ph.D. degree in engineering sciences from the Technical University of Lisbon. She is currently an Associate Professor with Nova School of Information Management, and the Head of SIMAQ Quality Management System. She has more than eighty refereed scientific publications, including a book, articles in international journals, and conference proceedings. She has lectured over a dozen short courses in statistics topics, as well as courses on SAS system programming. She collaborated in several research and development projects, both national and international, and she coordinated the GSIMCLI research project funded by FCT (Portuguese Science Foundation). Her research interest includes spatial statistics, particularly the modeling of spatial—temporal phenomena using geostatistics.

AQ:4

AQ:3





**Joaquín Torres-Sospedra** received the Ph.D. degree in ensembles of neural networks and machine learning from Universitat Jaume I in 2011. Since January 2020, he has been the Scientific Coordinator of UBIK Geospatial Solutions. He has authored more than 120 papers in journals and conferences. His current research interests include indoor positioning solutions-based on Wi-Fi & BLE, machine learning, and evaluation. He is the Chair of the IPIN International Standards Committee and IPIN off-site Competition.



**Marco Painho** received the degree in environmental engineering from the Universidade Nova de Lisboa, Portugal, the master's degree in regional planning from The University of Massachusetts, Amherst, and the Ph.D. degree in geography from the University of California, Santa Barbara, CA, USA. He is currently a Professor of Geographic Information Systems and Science with the Nova School of Information Management (NOVA IMS), Universidade Nova de Lisboa. He is the Coordinator of the Master's in Geographic Information Systems and Science and the Master of Science in Geospatial Technologies. He has over 30 years of experience in the GIS domain and coordinated over 100 projects in the application areas of the environment, natural resources management transportation, teaching among others. He is the author and editor of over 200 academic and professional publications.



**Joaquín Huerta** is currently a Full Professor with the Department of Information Systems, University Jaume I, Spain, where he teaches several courses related to GIS and Internet technologies. He is the Head of the GEOTEC Research Group, the Director of the Erasmus Mundus Master of Science in Geospatial Technologies degree program, run jointly with the universities of Münster and Nova de Lisboa. In addition to academic activities he is a founding member of UBIK Geospatial Solutions. He is and has been the principal investigator of several research projects, including EU projects as A-WEAR, GEO-C, and EUROGEOSS. His current research interests include indoor positioning, smart cities, mobile and web GIS applications, and augmented reality.

IEEE PROCEEDINGS