

# Data assimilation for sensing aided geolocation database

Jaakko Ojaniemi, Risto Wichman  
Aalto University, School of Electrical Engineering,  
Department of Signal Processing and Acoustics, Finland

**Abstract**—Cognitive radio systems aim to take advantage of the spatiotemporal empty spectrum without causing harmful interference towards the primary network by utilizing knowledge of the prevailing radio environment. The radio environment is typically modeled with propagation models or by interpolating spatially distributed field measurement data. This paper presents a practical online data assimilation method based on the ensemble Kalman filter for estimating the spatial correlation of the time-variant primary field strength from a collection of sensing samples. The correlation structure known as the variogram or covariance function is in turn used in the algorithms for radio environment mapping. Furthermore, it is shown that the proposed method provides significant reduction in the computation time compared to traditional sampling methods, thus, it offers an efficient real-time solution for state estimation in the future geolocation databases.

**Index Terms**—Cognitive radio, ensemble Kalman filter, sensing, radio environment mapping

## I. INTRODUCTION

The first realizations of cognitive radios are advancing towards commercial deployment as white space devices (WSD) capable of flexible spectrum utilization. White spaces are locally or temporally available frequencies that are principally occupied by primary use such as TV broadcasting. The currently prevalent view [1] to the utilization of such spectrum resources is that WSDs must query a geolocation database (GDB) through some dedicated channel to obtain information about available frequencies and related maximum transmission powers for the location of the WSD. The maximum throughput of white space network depends on the accuracy of the information provided by the GDB, as significant protection margins are required to minimize the worst-case interference in primary TV receivers.

A GDB is fundamentally based on field strength estimates for the primary service obtained using terrain based radio propagation models. However, the limited geographical information restricts the achievable accuracy of the field strength estimates. Radio environment mapping (REM) [2] has been introduced as an alternative or complementary procedure to radio propagation models. In concept of REM a database stores information of the estimated radio environment and uses this information for example to provide transmission power limits to WSDs. It is likely that at least a subset of WSDs will be capable of measuring and reporting the TV signal strength at their respective locations; such information can be used to improve the accuracy of a geolocation database.

Recently, there have been several studies on efficient algorithms for estimating the signal strength in unmeasured locations from a limited number of measurement samples. For example, inverse distance weighting based methods and kriging interpolation were compared in [3], and it was concluded that kriging is the most efficient estimator in terms of minimum mean squared error in situations where relatively large number of measurement points are available. A multivariate kriging method for incorporating field measurements into radio propagation models showed to provide further accuracy especially in undersampled areas [4].

Although kriging techniques offer the best estimation accuracy the modeling procedure requires knowledge of the spatial autocorrelation of the prediction variable. In cognitive radio systems, the state of the radio environment is susceptible to continuous fluctuations due to changes in propagation conditions or transmission parameters of the primary operator. To that end, active monitoring and reporting of the radio environment is important so that the prevailing correlation structure of the signal needed in the interpolation phase can be captured accurately.

Due to its recursive nature the Kalman filter (KF) provides a computationally practical method for estimating the state of a time-variant system corrupted by random noise. In the context of geolocation database system such conditions exist when the WSDs measure the dominant radio environment and send the measurement data to the GDB. It is reasonable to assume that such system contains huge amount of measurement data from multiple sources, and efficient on-line methods for data handling are needed. An approach was presented in [5] where the REM accepted only sensing data that contributed positively on the estimation accuracy in terms of minimizing the mean squared error. In [6] a fixed rank kriging method was introduced for operating with massive data sets.

As we are considering a dynamical system involving abrupt changes in the radio conditions more adaptive solutions are required. Consequently, we introduce an efficient real-time method for estimating the spatial correlation of the primary signal based on particle filter known as the ensemble Kalman filter, and show that the computation time can be decreased by several hundreds of percentages with negligible effect on the estimation accuracy. This makes the technique attractive for future use in GDBs that utilize spatial interpolation methods in radio environment mapping.

The paper is organized as follows: Section II introduces

the ensemble Kalman filter with a computational description. In Section III we present a method for linking the ensemble Kalman filter to estimate the spatial process from the sensing data. The results are presented in Section IV followed by the conclusion in Section V.

## II. ASSIMILATION OF SENSORY DATA

### A. Ensemble Kalman filter

Unlike the well-known KF which provides an optimal solution to the linear state estimation problem under Gaussian noise with known covariance, the ensemble Kalman filter (EnKF) is a suboptimal estimator for possibly non-linear dynamical systems.

In EnKF a Monte Carlo (MC) method is used to approximate the evolution of the state probability density described by the so called Fokker-Planck equation. Particularly, in the MC method an ensemble of the model states describes the current state space with the mean as the best estimate of the true state, and the spreading of the ensemble as the error variance. By integrating this ensemble in time through the non-linear process it is straightforward to approximate the necessary moments of the state probability density function. In contrast, in traditional KF the pdf of the process state is fully described by the mean and the covariance identified from the analytical examination. The MC method in EnKF is especially beneficial for systems of high order and huge amount of measurement data since the true covariance matrix is replaced with the sample covariance calculated from the selected ensemble. This is advantageous as the measurement and process covariances does not have to be known a priori.

For linear systems the EnKF converges to the solution of the KF as the number of ensemble members approaches infinity. For non-linear systems a widely used solution is the extended Kalman filter (EKF), which, however, involves calculation of the Jacobians for the non-linear process and measurement functions resulting in increased computational requirements. In addition, the EKF does not account the true non-linear dynamics as it linearizes about the current state and thus neglects some important statistical characteristics of the state probability density. Thus, it cannot be generalized as an optimal estimator.

Let the following non-linear system model describe the process state  $x$  at discrete time index  $k + 1$ :

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}_k \quad (1)$$

where  $f$  is a function which maps the current state to the the next state,  $\mathbf{x}_k, \mathbf{w}_k \in \mathbb{R}^n$ ,  $n$  being the order of the model,  $\mathbf{u}_k \in \mathbb{R}^m$ ,  $m$  being the dimension of the control input. The measurements are:

$$\mathbf{z}_k = h(\mathbf{x}_k) + \mathbf{v}_k \quad (2)$$

where  $h$  is a function which relates the state to the measurement,  $\mathbf{z}_k, \mathbf{v}_k \in \mathbb{R}^r$ , where  $r$  is the number of measurements. The Gaussian distributed random variables  $\mathbf{w}_k$  and  $\mathbf{v}_k$  are the uncorrelated process and measurement noise, respectively, with zero-mean and covariance matrices  $\mathbf{Q}_k$  and  $\mathbf{R}_k$ .

Both KF and EnKF aim to estimate the process state  $\mathbf{x}_k$  by minimizing the error between the estimated state and the true state using the measurements  $\mathbf{z}_k$ . Here we focus on the EnKF due to its computational efficiency. The filtering procedure is conceptualized as predict and assimilation (measurement update) phases.

1) *Prediction*: In EnKF, the approximate mean and the state error covariance matrix are calculated from the ensemble. First, let

$$[\mathbf{x}_k^{p,1}, \dots, \mathbf{x}_k^{p,N}]$$

represent the  $N$  ensemble members, where  $\mathbf{x}_k^p \in \mathbb{R}^n$  are the predicted state estimates.

The ensemble mean is

$$\bar{\mathbf{x}}_k^p = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_k^{p,i}, \quad (3)$$

where column vector  $\bar{\mathbf{x}}_k^p \in \mathbb{R}^n$ . The ensemble error matrix  $\hat{\mathbf{X}}$  can be calculated from the dispersion of the ensemble members around the mean

$$\hat{\mathbf{X}}_k^p = [\mathbf{x}_k^{p,1} - \bar{\mathbf{x}}_k^p, \dots, \mathbf{x}_k^{p,N} - \bar{\mathbf{x}}_k^p],$$

$\hat{\mathbf{X}}_k^p \in \mathbb{R}^{n \times N}$ . Similarly for the observation error matrix

$$\hat{\mathbf{Z}}_k = [\mathbf{z}_k^1 - \bar{\mathbf{z}}_k, \dots, \mathbf{z}_k^N - \bar{\mathbf{z}}_k]$$

with dimensions  $\mathbb{R}^{r \times N}$ . The prediction error covariance matrix for the ensemble is then

$$\mathbf{P}_{x_k} = \frac{1}{N-1} \hat{\mathbf{X}}_k^p (\hat{\mathbf{X}}_k^p)^\top, \quad (4)$$

and the measurement error covariance matrix

$$\mathbf{P}_{z_k} = \frac{1}{N-1} \hat{\mathbf{Z}}_k (\hat{\mathbf{Z}}_k)^\top. \quad (5)$$

2) *Assimilation*: The data assimilation is done by using the perturbed ensemble of the measurements. Since this ensemble is propagated through the system it allows to approximate the distribution of the states after undergoing a non-linear transformation. Given the measurements  $\mathbf{z}_k$  the starting point for data assimilation is to form an ensemble of perturbed measurements

$$\mathbf{z}_k^i = \mathbf{z}_k + \mathbf{v}_k^i, \quad i = 1 \dots N \quad (6)$$

where  $\mathbf{v}_k^i$  is a Gaussian noise term with zero mean and covariance  $\mathbf{R}_k$ . The data assimilation, or measurement update, is done as for KF but separately for each ensemble member

$$\mathbf{x}_k^{a,i} = \mathbf{x}_k^{p,i} + \hat{\mathbf{K}}_k (\mathbf{z}_k^i - h(\mathbf{x}_k^{p,i})), \quad (7)$$

where the Kalman gain is determined from the approximations of the error covariances using the observation matrix-free implementation as described in [8]

$$\hat{\mathbf{K}}_k = \mathbf{P}_{x_k} (\mathbf{P}_{x_k} + \mathbf{P}_{z_k})^{-1}. \quad (8)$$

The ensemble prediction is then performed using (1) by replacing  $\mathbf{x}_k$  with  $\mathbf{x}_k^{a,i}$  and  $\mathbf{w}_k$  with  $\mathbf{w}_k^{a,i}$ .

### B. Spatial correlation model

An essential concept in geostatistical modeling and interpolation is the variogram, which describes the spatial correlation as a function of distance. Generally, the correlation decreases with the distance, and the variogram is used to model the level and shape of the correlation. Interpolation methods such as kriging use the variogram for minimizing the interpolation error variances by determining distinct weights for the observation samples. Since the accuracy of the prediction, or interpolation error, fundamentally depends on the accuracy of the variogram it is important that the spatial correlation is modeled rigorously. The variogram can be described with [9]

$$\gamma(d) = \frac{1}{2N_z(d)} \sum_{(i,j) \in \mathcal{N}_{z_d}} |z_i - z_j|^2 \quad (9)$$

where  $\mathcal{N}_{z_d}$  is the set of pairs of observations  $i, j$  such that distance between measurement points  $|x_i - x_j| = d$ , where  $x$  represents the spatial coordinate,  $d$  is the distance class, and  $N_z(d)$  is the number of point pairs in that particular set.

The experimental variogram is fitted to a predefined variogram model, for example using the least squares method, to ensure that the estimation variance is positive and well-defined for all possible distances. There are several options for the variogram model, and the choice is typically made based on heuristics or on minimizing some error criterion. In our study the spherical model is used

$$\hat{\gamma}(d) = \begin{cases} b \left( \frac{3}{2} \frac{d}{a} - \frac{1}{2} \left( \frac{d}{a} \right)^3 \right), & d \leq a \\ b, & d > a \end{cases} \quad (10)$$

where  $b$  is the maximum value of the variogram function (sill variance),  $a$  describes the range or distance where the empirical variogram reaches its maximum, and  $d$  is the distance class.

### III. PROBLEM DEFINITION AND SIMULATION METHODS

Consider a cognitive radio system operating for example in TV white spaces. As described in Section I the performance of the system can be enhanced by updating the database with REM techniques by utilizing measurement samples from remote sensors such as WSDs. The data set consisting of the spectrum sensing samples will grow enormous in the course of time and thus requires efficient processing. Furthermore, the changes in the radio environment must be taken into account in the mapping procedure.

In this study, we demonstrate the use of ensemble Kalman filter for estimating the variogram calculated from the noisy sensing samples obtained in a changing radio environment. A realistic model of the network is constructed by using digital street data and sophisticated propagation modeling utilizing high-resolution digital terrain data. The sensing is then modeled by the predicted average field strength in given location corrupted by Gaussian noise. The simulation procedure is described in the following:

- 1) The radio environment is modeled by terrain based propagation model described in subsection III-A. This represents the true field strength from where the REM is estimated.
- 2) The changes in the radio environment are emulated by switching four transmitters on and off. The transmitters are located in the corners of network area shown in Fig.1.
- 3) The simulation consists of six phases. In the first four phases transmitters Tx1 to Tx4 are transmitting separately in each phase. In the fifth phase transmitters Tx1 and Tx4 are transmitting simultaneously. In the sixth phase all of the four transmitters are switched on simultaneously.
- 4) Each phase consists of 100 iterations. In each iteration a number of WSD sense the spectrum in the modeled area. The simulated sensing sample for a WSD is the value in its corresponding pixel in the propagation map (III-A) corrupted with zero mean random noise and standard deviation  $\sigma_{WSD}$ .
- 5) The locations for WSDs in each iteration are random but limited in the street network shown in Fig.1.
- 6) The database keeps a cumulative moving average of the field strength in the measured locations, that is, if a WSD senses an already measured pixel the GDB updates the average value with the new data.
- 7) The measurement data is uploaded to the GDB. The EnKF continuously estimates the variogram from the uploaded samples. This is described in subsection III-B.

#### A. Propagation model

The propagation prediction representing the true radio environment is implemented according to the guidelines presented in [10]. In addition to basic free-space propagation loss including short-term effects, the implemented model considers corrections due to different kinds of radio propagation phenomena. These include diffraction, tropospheric scatter, ducting and layer reflection/refraction, local clutter height, location variability, and building entry loss. The terrain data and parameters for the four prediction maps are the same except of the transmitter heights. The prediction map represents the Helsinki metropolitan area in Southern Finland spanning 12 km diagonally. The parameters are presented in Table I.

#### B. Mapping the spatial correlation

The variogram is estimated from the GDB in each iteration. Since the data set is growing in every iteration as more noisy samples are uploaded to the GDB, calculation of the variogram will become computationally demanding. However, the estimation simplifies by using EnKF.

In our model, the state of the system  $x_k$  is modeled as the values of the variogram functions calculated from a collection of predicted field strength values. As actual sensing data is migrated to the GDB the variogram functions will be calculated to form a set of measurements  $z_k$ . A major advantage in EnKF that it is not necessary to know the underlying process

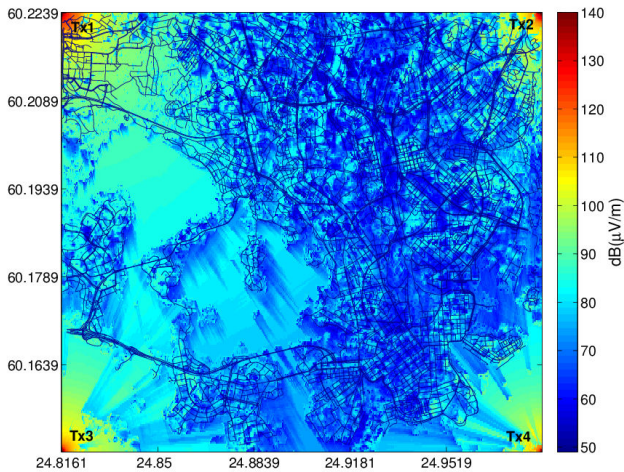


Fig. 1. Predicted field strength of four simultaneous transmissions, and the street network of the study area. The prediction map represents the Helsinki metropolitan area in Southern Finland, spanning 12 km diagonally.

Parameters for ITU-R P.1812	Value
ERP	1000 W
$h_{tx1} / h_{tx2} / h_{tx3} / h_{tx4}$	150m / 50m / 30m / 10m
$h_{rx}$	1.5m
Polarization	Horizontal
Frequency	490 MHz
Time percentage	50%
Location probability	50%
Resolution	30m $\times$ 30m
<b>Simulation parameters</b>	
$\sigma_{WSD}$	5.5dB
Nr. of sensors per iteration	30
Nr. of ensemble members	100
Nr. of distance classes	20

TABLE I  
SIMULATION PARAMETERS

since the state estimation is performed using the observations, where the true state is approximated by the ensemble mean. To simplify the modeling a one to one mapping between the process state and the observations is assumed, that is, each observation describes one state variable. Thus the dimension  $n$  of  $\mathbf{x}_k$  is equal to dimension  $r$  of  $\mathbf{z}_k$ .

The estimation procedure is described in the following. First, to initialize the filter and predict the process state,  $M$  locations are drawn randomly  $N$  times from the possible sensing locations in the geolocation database. Assuming that the location of the transmitter is known by the GDB, these initial predicted sample values at given distance from the transmitter are approximated as the free-space field strength corrupted with zero-mean Gaussian noise with standard deviation  $\sigma_{wsd}$ . The variogram is calculated for each of these  $N$  sets using (9) with  $r$  distance (lag) classes, corresponding to the number of the measurements in (2). The obtained set now forms the perturbed state predictions  $\hat{\mathbf{X}}_k^p$ .

Second, to obtain the measurements the variogram is calcu-

lated for the sensing data from the WSDs.  $N$  sets of the same  $M$  locations are sampled but the field strength is given by the procedure described in item 4) in the beginning of Section III. The variogram is calculated for each of these  $N$  sets. Similarly, the obtained set now forms the perturbed state measurements  $\hat{\mathbf{Z}}_k^p$ .

In turn, the database is approximating the mean process state by averaging (7) for each state to get  $\hat{x}_k^a$  by using (3)-(8) with the data described above. To complete the modeling, the variogram fitting is performed to  $\hat{x}_k^a$  using (10) after each predict and assimilation phase. The described procedure effectively maps the 2-dimensional spatial process to 1-dimensional function to be further used in geostatistical interpolation.

#### IV. RESULTS

As presented for example in [11] estimating the variogram reliably can require several thousands of samples depending on the spatial variation in the data. For total number of samples  $N_s$ , calculating the squared difference between point pairs in (9) requires  $\binom{N_s}{2}$  calculations. However, by distributing the variogram calculation for EnKF with  $N$  members with  $M$  samples per member the calculation reduces to  $N \binom{M}{2}$  operations. For example, computation for  $N_s = 10,000$  samples requires  $\sim 50M$  operations while distributing it to  $N = 100$  members with  $M = 100$  samples each requires  $\sim 0.5M$  operations. Using the latter method introduces additional computation in the order of  $\lfloor \Theta(r^2 N) \rfloor$  [8] for evaluating the Kalman gain in (8), however, resulting only in  $\sim 0.04M$  calculations with the parameters used.

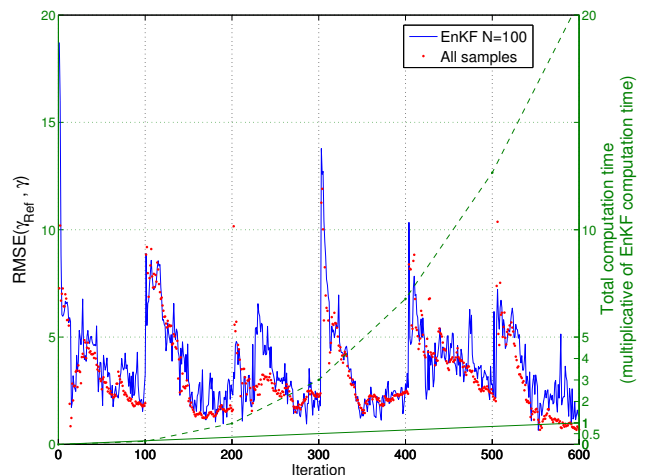


Fig. 2. Root mean square error between the reference variogram and the variogram calculated from the database content using the ensemble Kalman filter for sampling (blue), and all samples (red). The second ordinate (green) represents the duration for computing the variogram with EnKF (solid line) and comprehensively from all database samples (dashed line).

Nevertheless, there is a tradeoff between the estimation accuracy and computation time of the EnKF. This is shown in Fig. 2 as root mean square error (RMSE) between fitted variograms. The blue line represents the error for EnKF, and

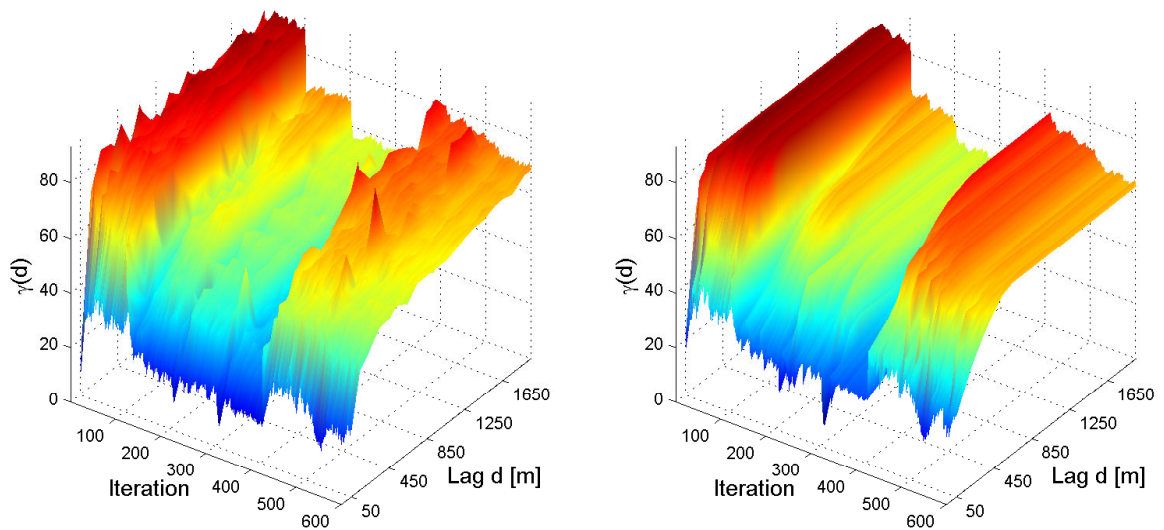


Fig. 3. Calculated variograms as a function of uploaded sensing samples (iteration). Left: Variogram calculated using EnKF. Right: Fitted variogram.

red dots represent the error for variogram calculated from all of the samples contained in GDB at the corresponding iteration. The reference variogram  $\gamma_{REF}$  for RMSE is calculated by sampling the model representing the real radio environment at 10,000 locations. However, as the figure indicates, the RMSE follows the ideal case (red dots) appropriately. In addition, the figure shows the computation time for calculating the variogram using the EnKF and entire content of the GDB. The EnKF is almost 20 times faster even in the short simulated case, thus, the tradeoff is reasonable. Note that the computation time for EnKF grows linearly since the number of samples per member in each iteration is fixed, while the computation time for comprehensive sampling increase with the factor  $\binom{N_s}{2}$ .

Figure 3 shows an example of the evolution of the variogram as new samples are arriving in the GDB. Clearly, the shape, range and maximum value of the variogram function follow the changing radio environment.

The results in terms of mean absolute error (MAE) after ordinary kriging (OK) [9] interpolation using the estimated variograms at first and last iteration in each phase are presented in Table II. The resulting radio environment map is obtained by considering 1000 spectrum samples in the network area and using 40 nearest samples per prediction location in the OK. The map is then compared against another set of 1000 samples from different locations in the original map. According to the results, MAE and the accuracy of the corresponding variogram improve significantly even in the relatively short simulation cycle consisting of 100 iterations.

## V. CONCLUSION

This paper presented a computationally practical method for data assimilation in geolocation database systems where abrupt changes in the radio environment cause reassessment

Variogram	P1	P2	P3	P4	P5	P6
$\gamma_{1st}$	7.40 dB	6.73 dB	5.59 dB	5.31 dB	6.17 dB	5.76 dB
$\gamma_{100th}$	5.05 dB	5.55 dB	4.75 dB	4.79 dB	4.57 dB	4.63 dB

TABLE II  
MEAN ABSOLUTE ERROR OF THE INTERPOLATED RESULTS USING THE ESTIMATED VARIOGRAMS FROM FIRST AND LAST ITERATION IN THE PHASES 1-6 OF THE SIMULATION.

of the parameters for the REM. Particularly, we modeled the 2-dimensional spatial process as 1-dimensional function, and linked the ensemble Kalman filter to estimate the mean state and covariance of the distance classes of the variogram. This procedure showed considerable improvement in the computational efficiency compared to traditional geostatistical sampling methods.

## ACKNOWLEDGEMENTS

This work was funded by Tekes, the Finnish Funding Agency for Technology and Innovation, in the WISE project [12] as the part of *Trial* technology program.

## REFERENCES

- [1] "Technical and operational requirements for the possible operation of cognitive radio systems in the white spaces of the frequency band 470-790 MHz," *ECC Report 159*, January 2011. Available online through <http://www.erodocdb.dk>
- [2] Zhao, Y., Morales, L., Gaeddert, J., Bae, K.K., Jung-Sun Um, Reed, J.H., "Applying Radio Environment Maps to Cognitive Wireless Regional Area Networks," *Proc. IEEE Symposium on New Frontiers in Dynamic Spectrum Access Networks (DySPAN 2007)*, Dublin, Ireland, April 2007.
- [3] Angeljicoski, M.; Atanasovski, V.; Gavrilovska, L., "Comparative analysis of spatial interpolation methods for creating radio environment maps," *19th Telecommunications Forum (TELFOR 2011)*, pp.334-337, 22-24 Nov. 2011

- [4] Ojaniemi, J., Kalliovaara, J., Poikonen, J., Wichman, R., "A practical method for combining multivariate data in radio environment mapping," Proc. *24th IEEE International Symposium on Personal Indoor and Mobile Radio Communications (PIMRC 2013)*, pp.729-733, London, UK, 8-11 Sept. 2013
- [5] Grimoud, S., Sayrac, B., Ben Jemaa, S., Moulines, E., "An algorithm for fast REM construction," Proc. *6th International ICST Conference on Cognitive Radio Oriented Wireless Networks and Communications (CROWNCOM) 2011*, pp.251-255, 1-3 June 2011
- [6] Riihijarvi, J.; Nasreddine, J.; Mahonen, P., "Demonstrating radio environment map construction from massive data sets," Proc. *IEEE International Symposium on Dynamic Spectrum Access Networks (DYSPAN 2012)*, pp.266-267, Bellevue, USA, 16-19 Oct., 2012
- [7] Evensen, G., "The Ensemble Kalman Filter: theoretical formulation and practical implementation," *Ocean Dynamics*, Vol. 53, No. 4, pp.343-367, 2003
- [8] Mandel, J., "Efficient Implementation of the Ensemble Kalman Filter," University of Colorado at Denver and Health Sciences Center, Tech. Rep. UCDHSC/CCM No. 231, May 2006
- [9] Wackernagel, H., *Multivariate Geostatistics*, Springer, Berlin, 1995.
- [10] "A path-specific propagation prediction method for point-to-area terrestrial services in the VHF and UHF bands," *ITU-R Recommendation P.1812-2*, Feb. 2012. Available: <http://http://www.itu.int/rec/R-REC-P.1812/en>
- [11] Zhang, H., Lan, Y., Lacey, R., Huang, Y., Hoffmann, W.C., Martin, B., Bora, C.G., "Analysis of variograms with various sample sizes from a multispectral image," *International Journal of Agricultural and Biological Engineering*, Vol. 2 No. 4, December, 2009
- [12] WISE project, <http://wise.turkuamk.fi>