

# Semi-supervised Learning for WLAN Positioning

Teemu Pulkkinen, Teemu Roos and Petri Myllymäki

Helsinki Institute for Information Technology HIIT  
PO Box 68, FI-00014 University of Helsinki, Finland  
`{firstname.lastname}@cs.helsinki.fi`

**Abstract.** Currently the most accurate WLAN positioning systems are based on the fingerprinting approach, where a “radio map” is constructed by modeling how the signal strength measurements vary according to the location. However, collecting a sufficient amount of location-tagged training data is a rather tedious and time consuming task, especially in indoor scenarios — the main application area of WLAN positioning — where GPS coverage is unavailable. To alleviate this problem, we present a semi-supervised manifold learning technique for building accurate radio maps from partially labeled data, where only a small portion of the signal strength measurements need to be tagged with the corresponding coordinates. The basic idea is to construct a non-linear projection that maps high-dimensional signal fingerprints onto a two-dimensional manifold, thereby dramatically reducing the need of location-tagged data. Our results from a deployment in a real-world experiment demonstrate the practical utility of the method.

**Keywords:** non-linear projection, manifold learning, wlan positioning, Isomap

## 1 Introduction

The need for special-purpose positioning systems for indoor use arises from the failure of established technologies, such as GPS, to properly locate and track objects in an indoor environment [8]. GPS signals tend to be weak when blocked by building walls, and even when a position is triangulated the accuracy is not sufficient for indoor use [4]. Several systems have been proposed that rely on the localized object carrying some kind of transceiver (RFID) [9] or infrared sensors built into the environment [14].

Recently, the interest in positioning based on wireless local area networks (WLANs), in particular, has grown significantly. This can be attributed to their wide use and distribution as well as the open standard which allows for requesting of signal strength information without separate authentication. WLAN-based systems have come a long way since the pioneering work of Bahl and Padmanabhan, who applied a nearest neighbor method on fingerprints composed of received signal strength indicator (RSSI) values [1]. Many of the most successful methods currently used in the field are probabilistic in nature ([6],[17], [23]). For a survey

on indoor positioning techniques, see [8]; for recent work, we refer the reader to [12],[22],[25].

Though WLAN fingerprinting approaches have achieved relatively good accuracy, and have found their way into some commercial services (e.g. [5]), the majority of location-based services are still based on GPS and other technologies [15]. One of the reasons to this is probably the manual effort required in *calibrating* fingerprinting-based methods: before the system can be used, fingerprints need to be recorded everywhere in the deployment area. Since the radio map created through this effort needs to be tied to real-world coordinates, it is also necessary to record the location of every fingerprint. This invariably requires human presence or other external location information (e.g., GPS, camera arrays) for the entirety of the calibration process.

We present a method for WLAN positioning wherein the fingerprinting approach is augmented with non-linear dimension reduction techniques. The main idea is to learn a low-dimensional, non-linear manifold that can represent the radio map, enabling better statistical modeling of the signal properties in complex multi-path environments. Once the manifold is constructed, we further propose a very simple method for mapping observation points attached to the manifold into geographical coordinates. Our approach is semi-supervised as the manifold learning phase is based on observing plain RSSI vectors without their geographical coordinates. A small sample of *key points* whose location is recorded are needed only to fix the mapping from the coordinate system of the manifold to geographical coordinates.

Earlier related work has focused on localization in sensor-networks. In the sensor-network localization problem a large set of sensor nodes communicate with other nodes in their proximity: Shang et al. [19] use the Isomap algorithm [21], and Patwari and Hero [13] use Laplacian eigenmaps to process binary connectivity data from each of the sensor nodes. Pan et al. [10, 11] apply Laplacian regularized least squares regression [2], without explicitly constructing a low-dimensional manifold; the drawback of this method is that the outcome is highly sensitive to the choice of the parameters controlling the regularization [24].

The rest of this paper is organized as follows. In Section 2, we lay out the basic concepts in semi-supervised learning, and in particular, manifold learning, including the specific non-linear approach (Isomap) used in this paper. In Section 3, we present the empirical framework and the details of the testing environment. Conclusions are summarized in Section 4.

## 2 The Semi-supervised Approach

Manifold learning methods attempt to find the defining features of a high-dimensional data set by reducing the dimensions (number of features) of the data to a more manageable level, usually two or three. The underlying assumption is that most of the variability in the data is concentrated on a low-dimensional (possibly non-linear) manifold embedded in the high-dimensional space. In our case, this is natural assuming that the signal characteristics are determined by

the location of the receiver, and that the dependency is smooth. If the possible locations are constrained to a flat two-dimensional surface, the resulting manifold is then two-dimensional as well. The crux of this approach is maintaining the pairwise distances between the fingerprints, at least locally, when they are mapped from the high-dimensional signal space to the low-dimensional manifold.

## 2.1 Isomap

One of the established manifold learning methods is the Isomap algorithm [21]. Isomap is based on the same principle as multidimensional scaling (MDS) in that, given a dissimilarity matrix, it tries to find a lower dimensional representation of the data such that the pairwise distances between the points are distorted as little as possible. One way to cast this as an optimization problem is to minimize the sum of squared deviations between the actual distances  $d_X(i, j)$ , and the distances in the new representation  $d_Y(i, j)$ :

$$\min_Y \sum_{i=1}^t \sum_{j=1}^t (d_X(i, j) - d_Y(i, j))^2, \quad (1)$$

If the original distances,  $d_X(i, j)$ , are Euclidean, MDS reduces to principal component analysis (PCA) [3]. Due to space limitations, we omit further details and refer the interested reader to [7].

Given a set of  $m$ -dimensional column vectors  $X = (\mathbf{x}'_1, \dots, \mathbf{x}'_n)$ , we denote by  $D = [d_X(i, j)]$  the matrix defined by their Euclidean distances. Further, we define  $B = HDH$ , where  $H$  is the symmetric centering matrix  $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$ , where  $\mathbf{1}$  denotes the all-ones column matrix, and  $\mathbf{1}^T$  its transpose. This implies that both the vector and column sums of  $B$  are null. Letting  $B = V\Lambda V^T$ , where  $\Lambda$  is a diagonal matrix, be the eigendecomposition of  $B$ , we obtain the eigenvectors as the columns of  $V$ , and the eigenvalues as the diagonal elements of  $\Lambda$ . The reconstruction obtained by using the  $l \geq 1$  largest eigenvalues,  $Y = V_p \Lambda_p^{\frac{1}{2}}$  is optimal in the sense of Eq. (1). An important observation is that if we replace the Euclidean distances  $d_X(i, j)$  by arbitrary dissimilarity values, which may or may not satisfy the properties of a valid distance metric, a solution can still be obtained by setting all negative eigenvalues (if any) to zero.

In the Isomap algorithm, the distances  $d_X(i, j)$  are obtained by constructing a neighborhood graph where each point  $\mathbf{x}_i$  is connected to its  $K$  nearest neighbors (in Euclidean distance). The length of an edge connecting two points is defined as their distance, and the distance  $d_X(i, j)$  between two points (that need not be neighbors) is then calculated as the sum of edge lengths along the shortest path connecting them. Applying the MDS algorithm as outlined above to the resulting distance matrix, yields a low-dimensional representation where the pairwise distances approximate path lengths along the neighborhood graph.

## 2.2 Manifold-Based Radio Map Learning

We now describe the application of Isomap in WLAN-based positioning. Consider a sample  $S = (\mathbf{s}'_1, \dots, \mathbf{s}'_n)$  of *fingerprints*, each of which is represented as a

vector  $\mathbf{s}_i = (s_{i1}, \dots, s_{ip})$  of RSSI values. The length of the vector,  $p$ , is defined as the number of access points (APs) in the WLAN network. The distance matrix  $X$  is then given by the Euclidean distance between the fingerprint vectors,  $d_X(i, j) = \|\mathbf{s}_i - \mathbf{s}_j\|_2$ . One of the practical problems that need to be solved is treating the occasionally unobserved RSSI values that show up as missing entries in the fingerprint vectors. Since the unobserved values are usually caused by too weak received signal, one reasonable solution is to replace all missing values by a small dummy value. In practice, we found that missing values typically result when the signal power drops below -100 dBm, and hence, we replaced all missing values by the constant -100 dBm, see [10, 17].

Another technical detail, albeit one that has a dramatic effect on the quality of the radio map produced by Isomap, is the choice of the neighborhood size,  $K$ . There is no universally good value, as appropriate values are determined by the variance of the observations perpendicular to the manifold relative to its curvature, and the sparseness of the available data [18, 20]. For too small a neighborhood, the neighborhood graph will not properly capture the geodesic distances on the manifold. Too large a neighborhood, on the other hand, risks creating “short circuits” that distort the topological properties of the manifold and make the algorithm unstable.

We propose to solve the neighborhood selection problem by exploiting additional information available in a set of fingerprints that are labelled by their geographical coordinates, which we call the *key points*. The method we propose below depends on being able to map points on the manifold onto a geographical coordinate system; we first describe a method for doing this.

### 2.3 Calibrating the Manifold to Geographical Coordinates

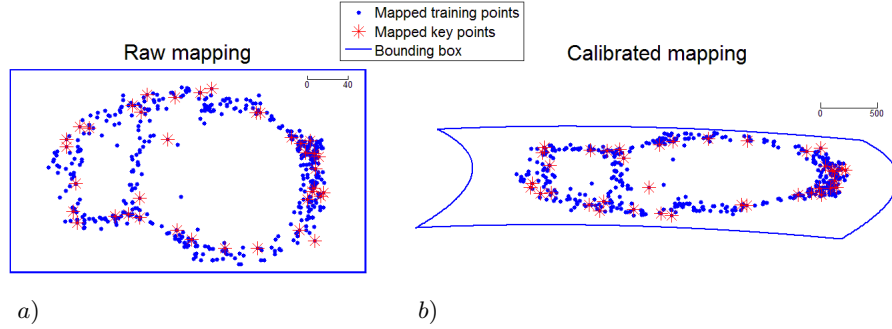
While the manifold learned by Isomap will reflect the topological structure of the area from which the data was collected, see Fig. 1a, it will usually not correctly match its metric properties such as lengths, angles, and curvature, which makes it unsuitable for positioning. This is corrected in what we call the *calibration* phase. We have found that the following very straightforward method is effective.

Assume that we have access to the precise location of  $n_{key}$  fingerprints, which we can without loss of generality assume to be the first  $n_{key}$  out of the total sample size of  $n$ . We denote the geographical coordinates of these *key points* by  $(g_i^{(x)}, g_i^{(y)})_{1 \leq i \leq n_{key}}$ . Denoting the manifold coordinates of the fingerprints by  $(m_i^{(x)}, m_i^{(y)})_{1 \leq i \leq n}$ , we map the manifold coordinates to geographical coordinates via

$$\begin{aligned} g_i^{(x)} &= \beta_x \tilde{\mathbf{m}}_i' + \epsilon_i^{(x)} \\ g_i^{(y)} &= \beta_y \tilde{\mathbf{m}}_i' + \epsilon_i^{(y)}, \end{aligned} \quad (2)$$

where  $\beta_x$  and  $\beta_y$  are both parameter vectors of length five, and

$$\tilde{\mathbf{m}}_i = (1, m_i^{(x)}, (m_i^{(x)})^2, m_i^{(y)}, (m_i^{(y)})^2), \quad (3)$$



**Fig. 1.** *a)* Manifold discovered by Isomap with the fingerprints on it, and *b)* the same manifold calibrated with geographical coordinates of a subset of key points (fingerprints marked with red stars).

are the regressor variables where we include the constant (intercept) term, both the manifold coordinates, as well as their squares. Note that the labeling of the manifold coordinates as  $x$  and  $y$  has no significance. The parameters  $\beta_x$  and  $\beta_y$  can be estimated by the standard least squares technique to minimize the sum of squares of the respective errors  $\epsilon_i^{(x)}$  and  $\epsilon_i^{(y)}$  for the key points  $1 \leq i \leq n_{key}$ . This provides an efficient way to map any point on the manifold, expressed as  $(m^{(x)}, m^{(y)})$  onto the corresponding geographical coordinates  $(g^{(x)}, g^{(y)})$ .

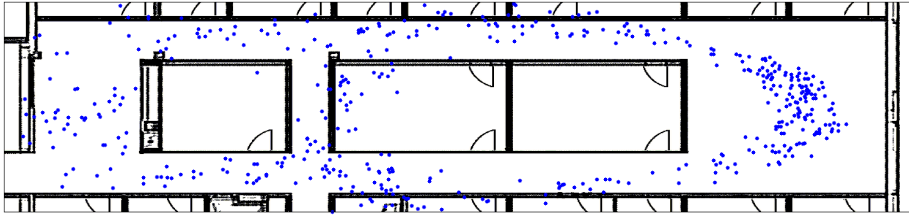
Figure 1 illustrates the process. The fact that the squares of the manifold coordinates are involved in Eq. (3) allows non-linear (namely quadratic) mappings, which is important since there is no guarantee that the correspondence between the learned manifold and the actual locations is linear. The non-linearity of the fit is clearly visible in the distortion of the bounding box in Fig. 1b. If a more generous set of key points is available, it may be useful to consider even more flexible mappings such as nonparametric regression, see [16].

Finally, we can use the error in the calibration mapping (2) to adjust the Isomap neighborhood size,  $K$ . We do this by trying different values between one and the total number of fingerprints (minus one), and choosing the one that leads to the mapping with the smallest error between the embedded key points and their actual (known) geographical coordinates:

$$\frac{1}{n_{key}} \sum_{i=1}^{n_{key}} \left( (g_i^{(x)} - \hat{\beta}_x \tilde{\mathbf{m}}_i')^2 + (g_i^{(y)} - \hat{\beta}_y \tilde{\mathbf{m}}_i')^2 \right).$$

## 2.4 Positioning

Positioning new fingerprints is relatively straightforward once the manifold has been learned and calibrated with the key points. There are various ways to map new fingerprints onto the manifold, and thence to geographic coordinates. We



**Fig. 2.** Plot of embedded fingerprints

choose to use the *k-nearest neighbors* method, selecting the  $k$  nearest fingerprints (not necessarily any of the key points), and then letting the manifold coordinates of the new fingerprint be given by the average of the coordinates of the selected  $k$  fingerprints. The latter are directly obtained from Isomap output. The resulting manifold coordinates are then mapped to geographical coordinates, providing the position estimate, by Eq. (2). A comparison of alternative positioning methods in combination with manifold approaches is an interesting topic for future work.

### 3 Deployment and Results

We deployed the system in a real-world office building at the Department of Computer Science, University of Helsinki. The deployment area covered hallways and an adjoining open space used as a meeting space. The total area of the environment was about  $24 \text{ m} \times 7 \text{ m}$ . The data recording, processing, and most positioning tests were performed with a Samsung NC10 Netbook, running Ubuntu Linux 9.10, equipped with an Atheros AR5007EG Wireless network adapter, complying to the 802.11b/g standard. The total number of fingerprints used for learning and calibrating the manifold was  $n = 437$ , of which  $n_{key} = 38$  were used as key points. We reserved an additional  $n_{test} = 66$  points for testing purposes.

The Isomap neighborhood size that was found to minimize the error in mapping the key points was 15. This left the average error of 1.9 m. Among the 66 test fingerprints collected separately, the mean positioning error was 2.0 m, and the median error was 1.5 m. Plotting the calibrated points onto the floor plan, we can clearly see the shape of the hallway in the mass of points, see Fig. 2. A majority of the points mapped to the hallway respect the infrastructure. It is clear that the hallways insulate the WLAN signal and create unique signatures. The mapping of fingerprints in the open space was not as distinct, however. This was most likely caused by the lack of attenuating infrastructure, making it hard to distinguish between the fingerprints from different ends of the space.

We have also carried out experiments in other environments with somewhat varying results; details are omitted due to space restrictions. Future research will benefit from an investigation of the factors most affecting the outcome.

## 4 Conclusion

We presented a WLAN positioning approach where high-dimensional signal fingerprints are represented as points on a two-dimensional manifold. For the manifold learning phase, we used the Isomap algorithm. Our contributions include a straightforward method for mapping points on the Isomap manifold to a geographical coordinate system by taking advantage of a relatively small subset of the fingerprints whose precise location is known. This also allowed us to choose the neighborhood size, a central (and only) parameter in Isomap, in a principled way by minimizing the error in the resulting coordinate mapping.

The main benefits of our method are: more robust estimation of the RSSI variability due to the lower dimensionality of the estimated model, and even more importantly, reduction in the effort required to collect measurement data. The latter feature boosts the cost-effectiveness of the fingerprinting approach both in terms of initial set-up as well as maintenance, which may finally enable WLAN-based indoor positioning to become the method of choice for future location-based services. Exploring the exact tradeoff between the number of labelled examples (and thus the deployment cost), and accuracy is a most urgent topic for investigation, which, however, is beyond the scope of this paper.

**Acknowledgments.** This work was supported in part by the European Commission under the PASCAL Network of Excellence.

## References

1. Bahl, P. and Padmanabhan, V. N.: RADAR: An In-building RF-based User Location and Tracking system. In: 19th Conference of IEEE Computer and Communications Societies, pp 775–784. IEEE Computer Society, Piscataway (2000)
2. Belkin M., Niyogi P., and Sindhvani, V.: On Manifold Regularization. In: 10th International Workshop on Artificial Intelligence and Statistics, pp. 17–24. IOS Press, Amsterdam (2005).
3. Cox, T. F. and Cox, M. A. A.: Multidimensional Scaling. Chapman & Hall, London (2001)
4. Dedes, G. and Dempster, A.: Indoor GPS: Positioning Challenges and Opportunities. In: 62nd Vehicular Technology Conference, pp. 412–415. IEEE Press, Piscataway (2005).
5. Ekahau, Inc. RTLS, <http://www.ekahau.com>
6. Ferris, B., Hahnel, D. and Fox, D.: Gaussian Processes for Signal Strength-Based Location Estimation. In: Robotics: Science and Systems, pp. 1–8. MIT Press, Cambridge (2006)
7. Ghodsi, A: Dimensionality Reduction – A Short Tutorial. Technical report, University of Waterloo (2006)
8. Liu, H., Darabi, H., Banerjee, P. and Liu, J.: Survey of Wireless Indoor Positioning Techniques and Systems. In: IEEE T. Syst. Man. Cyb. 37, 1067–1080 (2007)
9. Ni, L. M., Liu, Y., Lau, Y. C. and Patil, A. P.: LANDMARC: Indoor Location Sensing Using Active RFID. In: Wireless Networks. 10, 701–710 (2004)

10. Pan, J. J., Yang, Q., Chang, H. and Yeung, D.-Y.: A Manifold Regularization Approach to Calibration Reduction for Sensor-Network Based Tracking. In: 21st National Conference on Artificial Intelligence, pp. 988–993. AAAI Press, Menlo Park (2006).
11. Pan, J. J. and Yang, Q.: Co-localization from Labeled and Unlabeled Data Using Graph Laplacian. In: 20th International Joint Conference on Artificial Intelligence, pp. 2166–2171. Morgan Kaufmann, San Francisco (2007)
12. Papapostolou, A. and Chaouchi, H.: WIFE: Wireless Indoor Positioning Based on Fingerprint Evaluation. In: Fratta, L., Schulzrinne, Henning, Takahashi, Y. and Spaniol, O. (eds) NETWORKING 2009. LNCS, vol. 5550, pp. 234–247. Springer, Heidelberg (2009)
13. Patwari, N. and Hero, A. O.: Adaptive Neighborhoods for Manifold Learning-based Sensor Localization. In: IEEE 6th Workshop on Signal Processing Advances in Wireless Communications, pp. 1098–1102. IEEE Press, Piscataway (2005).
14. Petrellis, N., Konofaos, N. and Alexiou, G.: A Wireless Infrared Sensor Network for the Estimation of the Position and Orientation of a Moving Target. In: Third International Mobile Multimedia Communications Conference, pp. 1–4. ICST, Brussels (2007)
15. Raper, J., Gartner, G., Karimi, H. and Rizos, C.: Applications of Location-based Services: A Selected Review. *J. Location-based Services* 1, 89–111 (2007).
16. Rasmussen, C. E. and Williams, C. K. I.: *Gaussian Processes for Machine Learning*. MIT Press, Cambridge (2005)
17. Roos, T., Myllymäki, P., Tirri, H., Misikangas, P. and Sievänen, J.: A Probabilistic Approach to WLAN User Location Estimation. *Int. J. Wireless Information Networks*. 9, 155–164 (2002)
18. Samko, O., Marshall, A. D. and Rosin, P. L.: Selection of the Optimal Parameter Value for the Isomap Algorithm. *Pattern Recogn. Letters*. 27, 968–979 (2006)
19. Shang, Y., Ruml, W., Zhang, Y. and Fromherz, M. P. J.: Localization from Mere Connectivity. In: 4th ACM International Symposium on Mobile Ad Hoc Networking & Computing, pp. 201–212. ACM Press, New York (2003)
20. Shao, C., Huang, H. and Wan, C.: Selection of the Suitable Neighborhood Size for the Isomap Algorithm. In: International Joint Conference on Neural Networks, pp. 300–305. IEEE Press, Piscataway (2007).
21. Tenenbaum, J. B., de Silva, V. and Langford, J. C.: A Global Geometric Framework for Nonlinear Dimensionality Reduction. *Science*. 290, 2319–2323 (2000)
22. Yeung, W. H., Zhou, J.-Y., and Ng, J. K.: Enhanced Fingerprint-Based Location Estimation System in Wireless LAN Environment. In: Denko, M. K., Shih, C., Li, K.-C., Tsao, S.-L., Zeng, Q.-A., Park, S. H., Ko, Y.-B, Hung, S.-H. and Park, J. H. (eds) EUC 2007. LNCS, vol. 4809, pp. 273–284. Springer, Heidelberg (2007)
23. Youssef, M., Agrawala, A. and Shankar, A. U.: WLAN Location Determination via Clustering and Probability Distributions. In: 1st IEEE International Conference on Pervasive Computing and Communications, pp. 1–8. IEEE Press, Piscataway (2003)
24. Yuan, J., Li, Y., Liu, C., Zha, X. F.: Leave-One-Out Cross-Validation Based Model Selection for Manifold Regularization. In: Yuan, J., Li, Y.-M., Liu, C.-L. and Zha, X.F. (eds) ISNN 2010. LNCS, vol. 6064, pp. 457–464. Springer, Heidelberg (2010)
25. Zhang, M. and Zhang, S.: An Accurate and Fast WLAN User Location Estimation Method Based on Received Signal Strength. In: Zhang, M. and Zhang, S. (eds) ICCS 2007. LNCS, vol. 4489, 58–65. Springer, Heidelberg (2007)