

In this section, we describe uncertainty sampling for kNN algorithm newly designed for kNN classifier. This is extended version of previous uncertainty sampling for nearest neighbor (NN) referred to [8].

Uncertainty Sampling for NN (US_NN) algorithm calculates the uncertainty of samples and preferentially selects samples with high uncertainty. Uncertainty can be calculated based on least confident, margin, or entropy metric [5]. However, these uncertainty metrics are difficult to apply to the nearest neighbor (NN) classifier, which is a deterministic classifier since they are based on stochastic modeling techniques. The uncertainty metric for NN classifier has been defined in [7] as followings.

$$x_H^* = \underset{x}{\operatorname{argmax}} \left[- \sum_{0 \leq i < |L|} P_\theta(y_i | x) \log P_\theta(y_i | x) \right],$$

$$P(y|x) = \frac{a_y}{\sum_{0 \leq i < |Y|} d_{y_i}}, \quad (2)$$

$$d_y = \min_{s \in S} (\delta(s, x)).$$

where Y is a set of all classes, δ is a distance algorithm, and S is a set of all samples of class y . The function (2) is an entropy-based uncertainty function and does not use Bayesian probabilistic likelihood, but use distance-based likelihood. Likelihood $P(y/x)$ in function (2) means ratio of distance between a sample x and corresponding nearest neighbor in class y to sum of distances d_{y_i} ($0 \leq i < |Y|$). The d_{y_i} is a distance between a sample x and nearest neighbor in class y_i .

Uncertainty Sampling for kNN (US_kNN) algorithm is similar to the US_NN algorithm except that it considers k nearest neighbors. Since the US_NN algorithm considers only closest samples for each class, the effect of outliers can be significant. We designed US_k-NN algorithm to overcome the problem of US_NN algorithm.

The US_k-NN algorithm is basically similar to the function (2). The difference is that US_kNN algorithm considers all class candidates that can be classified as class y in kNN classifier to calculate the distance d_y . The kNN classifier classifies the sample x as a class y , if the majority of the labels of k closest samples with x are y . The function $I(Y, k, y)$ in (3) outputs a set of all class candidates, Y_C .

$$I(Y, k, y) = \{(i_0, \dots, i_{k-1}) \mid i_{0, \dots, \lfloor \frac{k}{2} \rfloor - 1} = y, i_{\lfloor \frac{k}{2} \rfloor, \dots, k-1} = I_H(Y, \bar{k})\} \quad (3)$$

where Y is a set of all classes, $\bar{k} = 1 - \lfloor \frac{k}{2} \rfloor$, and $I_H(Y, \bar{k})$ is an extractor that outputs $|Y|^{\bar{k}}$ combinations with repetition from a set $Y - \{y\}$.

Now, we can calculate d_y through the Distance of kNN Candidates (DKC) algorithm including function (3). The d_y is the average of d_{kC} which is the average of k distances $\{d_{kC}(y_0), \dots, d_{kC}(y_{k-1})\}$. For each distance $d_{kC}(y_i)$, it represents the distance between the sample x and the closest sample x' with one of the label y_i in Y_C . Because one or more labels in Y_C are the same, we need the function (4) that is able to compute the distance between the sample x and the α -th closest sample with label y .

$$kMinDist(x, y, \alpha) = \min_{s \in S} (\delta(s, x)) \quad (4)$$

where δ is a distance algorithm, S is a set of all samples in a class y , and $\min_{\alpha}(\cdot)$ is the function that outputs α -th minimum value.

After getting d_{kC} with the number of $|Y_C|$, we can get d_y by averaging out the all d_{kC} for kNN candidates.

[DKC algorithm]
 Input : (x, y, k, Y)
 Output : d_y
 Algo:
 1. Y_c ← I(Y, k, y)
 2. sum_y = 0
 3. For {y₀, ..., y_{k-1}} in Y_c:
 4. For i = 0 to k-1:
 5. C(y_i) = 0 //initialization
 6. sum_{kC} = 0
 7. For y_i in {y₀, ..., y_{k-1}}:
 8. sum_{kC} = sum_{kC} + kMinDist(x, y_i, C(y_i))
 9. C(y_i) = C(y_i) + 1
 10. d_{kC} = sum_{kC} / k
 11. sum_y = sum_y + d_{kC}
 12. d_y = d_{kC} / |Y_c|
 13. return d_y

V. EMPIRICAL EVALUATION

In this section, we evaluate the performance of the proposed AR-WFL system using the dataset described in section 2. Two update algorithms are evaluated in detail; Naive Incremental Learning (NIL) and Active Learning with uncertainty sampling for kNN (AL_US_kNN). We used the first dataset, DS0829, as a Wi-Fi fingerprint database and others are used for test datasets. Classification algorithm is kNN and scikit-learn library in python is used for evaluation. We evaluate the performance by comparing accuracy of room-level classification.

$$accuracy = \frac{\sum_{r \in R} correct(S_r)}{|S|} \quad (5)$$

where S is all test samples, S_r is a subset of S whose samples are labeled with a room r , and $correct(S_r)$ is a subset of S_r whose samples are predicted as a room r .

The proposed continuous active learning system is a crowdsourcing based system and aims to improve the performance while minimizing the number of active queries. Thus, we checked the change of performance varying the number of samples selected through the selective sampling algorithms. The number of samples selected from the selective sampling algorithm is varied from 10 to 200 in units of 10. We limited the number of samples to 200 because the minimum number of samples in our collected datasets is 200.

Figure 5 shows the positioning accuracy according to the number of samples of active query for AL_US_kNN algorithm. The y-axis of the chart represents the average accuracy for the all test data from DS0831 to DS1028. There are 3 highest accuracies in the chart; (60, 94.38%), (170, 94.41%), and (200,

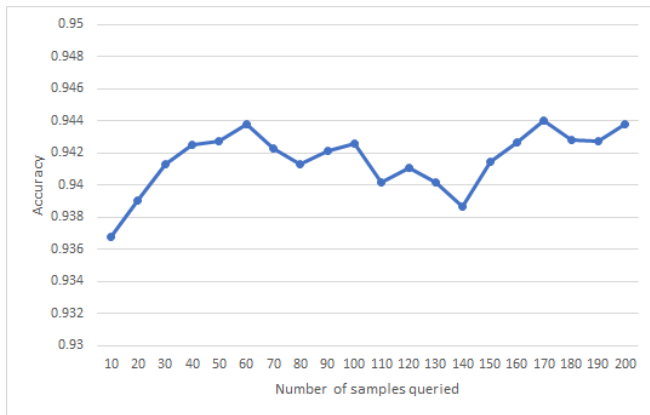


Fig. 5. Positioning accuracy of AL_US_kNN according to the number of active queries.

94.38%). Although maximum accuracy is when the number of samples queried is 170, we must consider the usability which means user intervention. The lower queries, the better usability so it would be more reasonable that we choose the number of queries as 60.

For performance comparison, we also evaluated the NIL and the case that update module is not be applied called NoUpdate. As shown in figure 6, for NoUpdate, the maximum accuracy is 97.53% for DS0919 and the average accuracy is 92.55%. For NIL, the maximum accuracy is 99.04% for DS0929 and the average accuracy is 93.30%. For AL_US_kNN, the maximum accuracy is 98.53% for DS0929 and the average accuracy is 94.38%.

NIL, the easiest way to think about database updates, is a fully automated method that does not require user participation. When the NIL technique was applied, the average accuracy increased by 0.75%p compared to NoUpdate. This shows that the NIL method can mitigate the performance degradation due to aging of the Wi-Fi fingerprint database.

Active learning method, AL_US_kNN, achieved the higher average accuracy than NIL. When the AL_US_kNN technique was applied, the average accuracy increased by 1.83%p compared to NoUpdate when the number of samples queried 60. It demonstrates that crowdsourcing based AR-WFL system using selective sampling algorithm can cope with the aging problem more efficiently in terms of accuracy than NIL method.

VI. RELATED WORKS

Wi-Fi fingerprint-based indoor localization technology has been continuously studied for about 20 years, starting with RADAR [2]. RADAR is the first supervised learning based WFL system that collects signal strength from 3 base stations, builds WF database, and estimates the location with kNN classifier.

Since RADAR, various WFL systems have been studied. According to [3], the study of WFL systems can be divided into two categories, one for high accuracy and the other for efficient deployment. For higher accuracy, [9] analyzed the wireless channel characteristics with temporal variation and spatial variation, [10-13] uses feature selection methods based on constancy, strength, or coverage of public Wi-Fi APs, and [14]

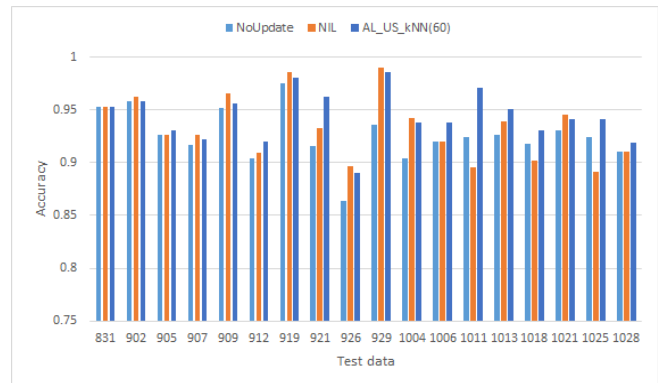


Fig. 6. Comparison of positioning accuracy for each test data with NoUpdate, NIL, and AL_US_kNN. The number of samples actively queried of AL_US_kNN is 60.

analyzed probability distribution of Wi-Fi fingerprint database called IGDG (Improved Double-peak Gaussian Distribution).

For efficient deployment, researches which are to reduce the site-survey for Wi-Fi fingerprint database [15-19], and keeping WF database up-to-date [20-23] have been conducted. Wi-Fi fingerprint database is constructed and updated based on the data collected by a person who walks around the building. This causes an increase of labor costs and a decrease of efficiency to build and manage the WFL system. Therefore, these studies aim to design an automated system that minimizes human intervention.

For keeping database up-to-date, [20] proposed LeManCoR system based on manifold co-regularization, [21] proposed the WFL system with dead reckoning that can minimize the reference points to be labeled by a human annotator. [22] proposed a crowdsourcing-based system, which assumes that users have a navigation solution called T-PN. T-PN operates based on inertial sensors and provides LLH (latitude, longitude, height) information. The system automatically updates the Wi-Fi fingerprint database if it detects the difference between user's current Wi-Fi fingerprint and that of the database although LLH is the same within thresholds. This system has the advantage of enabling database update without user's active participation, but it could be a limitation that system depends on T-PN solution. [23] proposed the automated system to update Wi-Fi fingerprint database by continuously tracking changes in signal strength using wireless sensor network (WSN). This system can automatically collect the labeled Wi-Fi fingerprint data based on WSN which sensors are pre-installed in fixed locations. Thus it has the advantage of instantly reflecting the change of Wi-Fi environment but disadvantage of being required the additional infra, WSN.

The active learning technique used in this paper has been studied in the field of machine learning. According to [4], active learning techniques have been studied for probabilistic classifiers and non-probabilistic classifiers. In particular, [7] proposed a selective sampling algorithm for NN classifier, which is a non-probabilistic classifier.

Active learning techniques are designed to solve the problem that many unlabeled data exist but obtaining labels of them are expensive. Since the WFL system is difficult to acquire labeled data, it is considered that the active learning scheme is applicable to the WFL system. There are no previous works to

apply the active learning scheme to the WFL system to the best of our knowledge.

VII. CONCLUSION

In this paper, we proposed AR-WFL system with continuous active learning to solve the aging problem in WFL system. The AR-WFL system is a crowdsourcing-based system that does not assume a dedicated annotator. We applied selective sampling algorithms to minimize user participation for updating Wi-Fi fingerprint database. For the performance evaluation of the AR-WFL system, 19 datasets were collected for 2 months in a target building with 8 rooms. Based on the collected data, uncertainty sampling for k-NN algorithm were implemented to compare the accuracy of location estimation. As a result, it showed that average accuracy is increase by 1.83%p compared with NoUpdate.

ACKNOWLEDGMENT

This work was supported by Institute for Information & communications Technology Promotion(IITP) grant funded by the Korea government(MSIT) (No.2015-0-00168, Development of Universal Authentication Platform Technology with Context-Aware Multi-Factor Authentication and Digital Signature)

REFERENCES

- [1] U.S. Environmental Protection Agency Green Building Workgroup, "Buildings and their impact on the environment: A statistical summary," 2009.
- [2] Paramvir Bahl and Venkata N. Padmanabhan, "RADAR: An In-Building RF-based User Location and Tracking System," *Tech. Rep. MSR-TR-00-12*, Microsoft Research, Feb. 2000.
- [3] Suining He and S.-H. Gary Chan, "Wi-Fi Fingerprint-Based Indoor Positioning: Recent Advances and Comparisons," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 466-490, First quarter 2016.
- [4] Burr Settles, "Active Learning Literature Survey," *Computer Sciences Technical Report 1648*, University of Wisconsin-Madison, Jan. 2010.
- [5] Burr Settles, "An Analysis of Active Learning Strategies for Sequence Labeling Tasks," *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp.1070-1079, Oct. 2009.
- [6] Prateek Jain and Ashish Kapoor, "Active Learning for Large Multi-class Problems," *IEEE Conference on Computer Vision and Pattern Recognition*, pp.762-769, 2009.
- [7] Michael Lindenbaum, Shaul Markovitch, and Dmitry Rusakov, "Selective Sampling for Nearest Neighbor Classifiers," *Machine Learning*, 54, pp.125-152, 2004.
- [8] Martina Hasenjaeger and Helge Ritter, "Active Learning with Local Models," *Natural Processing Letters*, pp.107-117, 1998.
- [9] Moustafa Youssef and Ashok Agrawala, "The Horus WLAN Location Determination System," *MobiSys '05*, pp.205-218, June 2005.
- [10] Arsham Farshad, Jiwei Li, Mahesh K. Marina, and Francisco J. Garcia, "A Microscopic Look at WiFi Fingerprinting for Indoor Mobile Phone Localization in Diverse Environments," *2013 International Conference on Indoor Positioning and Indoor Navigation (IPIN2013)*, Oct. 2013.
- [11] Jungmin So, Joo-Yub Lee, Cheal-Hwan Yoon, and Hyunjae Park, "An Improved Location Estimation Method for Wifi Fingerprint-based Indoor Localization," *International Journal of Software Engineering and Its Applications*, vol. 7, no. 3, pp. 77-86, May 2013.
- [12] Pei Jiang, Yunzhou Zhang, Wenyan Fu, Huiyu Liu, and Xiaolin Su, "Indoor Mobile Localization Based on Wi-Fi Fingerprint's Important Access Point," *International Journal of Distributed Sensor Networks*, vol. 2015, pp. 1-8, 2015.

- [13] Youngsam Kim and Soohyung Kim, "Rethinking of Feature Selection Methods for Room-Level Localization Using Public APs," *19th International Conference on Advanced Communication Technology (ICTACT2017)*, pp. 24-28, Feb. 2017.
- [14] Lina Chen, Binghao Li, Kai Zhao, Chris Rizos, and Zhengqi Zheng, "An Improved Algorithm to Generate a Wi-Fi Fingerprint Database for Indoor Positioning," *Sensors2013*, 13, pp. 11085-11096, Aug. 2013.
- [15] Anshul Rai, Krishna Kant Chintalapudi, Venkata N. Padmanabhan, and Rijurekha Sen, "Zee: Zero-Effort Crowdsourcing for Indoor Localization," *MobiCom '12*, pp. 293-304, Aug. 2012.
- [16] He Wang, Souvik Sen, Ahmed Elgohary, Moustafa Farid, Moustafa Youssef, and Romit Roy Choudhury, "No Need to War-Drive: Unsupervised Indoor Localization," *MobiSys '12*, pp. 197-210, June 2012.
- [17] Ahmad Abadleh, Sangyup Han, Soon J. Hyun, Ben Lee, and Myungchul Kim, "ILPS: Indoor Localization using Physical Maps and Smartphones Sensors," *15th International Symposium on a World of Wireless, Mobile and Multimedia Networks (WoWMoM)*, June 2014.
- [18] Chensu Wu, Zheng Yang, Yunhao Liu, and Wei Xi, "WILL: Wireless Indoor Localization without Site Survey," *IEEE Transactions on Parallel and Distributed Systems*, vol. 24, no. 4, pp. 839-848, Apr. 2013.
- [19] Chensu Wu, Zheng Yang, and Yunhao Liu, "Smartphones Based Crowdsourcing for Indoor Localization," *IEEE Transactions on Mobile Computing*, vol.14, no.2, pp. 444-457, Feb. 2015.
- [20] Sinno Jialin Pan, James T. Kwok, Qiang Yang, and Jeffrey Junfeng Pan, "Adaptive Localization in a Dynamic WiFi Environment Through Multi-view Learning," *AAAI'07*, vol. 2, pp. 1108-1113, July 2007.
- [21] Le T. Nguyen and Joy Zhang, "Wi-Fi fingerprinting through Active Learning using Smartphones," *UbiComp '13*, pp.969-976, Sep. 2013.
- [22] Yuan Zhuang, Zainab Syed, You Li, and Naser El-Sheimy, "Evaluation of Two WiFi Positioning Systems Based on Autonomous Crowdsourcing of Handheld Devices for Indoor Navigation," *IEEE Transactions on Mobile Computing*, vol. 15, no. 8, pp. 1982-1995, Aug. 2016.
- [23] Walter Balzano, Aniello Murano, and Fabio Vitale, "WiFACT - Wireless Fingerprinting Automated Continuous Training," *30th International Conference on Advanced Information Networking and Applications Workshops(WANIA)*, pp. 75-80, Mar. 2016.



learning, and security protocols.

Youngsam Kim received a B.S. (2009) degree in computer engineering from Chungbuk National University and a M.S. (2011) in information security engineering from the University of Science and Technology in South Korea. In 2011, he joined future internet research team in NIMS as a researcher. Currently, he is a researcher at the Electronics and Telecommunications Research Institute. His research interests include context-aware authentication, machine



biometrics, identity management, network and system security.

Soohyung Kim received the B.S. and M.S. degrees in computer science from Yonsei University, Seoul, Korea, in 1996 and 1998. He received the Ph.D. degree in computer science from Korea Advanced Institute of Science and Technology (KAIST), Daejeon, Korea in 2016. He is currently a Director of Information Security Research Division in Electronics and Telecommunications Research Institute (ETRI), Daejeon, Korea. His research interests include payment system,