# Quality Evaluation of Crowdsensed Fingerprints for Indoor Localization

*Abstract*—In crowdsensing, the dataprocurement is crucial

## I. INTRODUCTION

### A. Indoor localization

The past decade has witnessed a flourishing of indoor localization systems based on wireless techniques [?], where the fingerprinting based methodology has been widely adopted due to its convenient deployability [?], [?]. The fingerprinting based indoor localization system has two phases: In the offline phase, the site surveyor observes the received signal strength (RSS) of Wi-Fi access points (APs) termed as RSS fingerprints at each reference point, and submit the fingerprints and the location information of the reference point to the localization database; in the online phase, a user needs localization service could submit the observed fingerprints to the database, which then returns the location of the reference point that matches the fingerprints best as the estimated location of the user.

The fingerprinting based method utilizes Wi-Fi APs widely existing in buildings and has no need for other dedicated infrastructure; however, the site survey in the offline phase requires substantial efforts, which is hardly accomplished by any single entity. The recent advances of fingerprinting localization systems utilize mobile crowdsensing approach to collect fingerprints [?], [?], [?], [?], [?], [?]. Mobile crowdsensing is a cost-effective approach to collect large scale data for mobile applications, where individuals with hand-held mobile devices collectively contribute sensing data so that information of certain events could be retrieved [?], [?]. Although sensing participants could receive certain rewards for the efforts and resources spent on the sensing activity, the cost of mobile crowdsensing is still much lower than deploying the dedicated sensing networks [?].

As the crowdsensing data are collected by unprofessional participants with non-dedicated equipment, the sensing data obtained are usually with considerable noise. The quality of the sensing data is the crux for evaluating contribution of the participants, which is the vitally important for effective utilizing rewards to incentivize participants to accomplish sensing tasks satisfactorily. However, how to evaluate the quality of the crowdsensing data is a challenging issue, because there is no ground truth for the collected data to be compared with. Efforts have been made to evaluate the crowdsensing data quality [?], and the task allocation scheme [?], [?]. and incentive mechanisms considering the data quality are proposed.

While the efforts have been made for the evaluation of the quality of the data, the RSS data procurement still remains not fully studied yet. Some quality-driven incentive mechanism[?], [?], [?], [?], [?] were proposed, however, the state of the art method for crowdsensing data collection still focus on the incentive of workers and the utility of the platform. The economical problem is considered in [?], however, the budget constraint of the platform and is not included. Besides, all the work listed above do not consider the situation when data is coming in a sequential order and only available in each round. How to acquire the high-quality data that is in a sequential order given the limited budget still requires more thorough investigation, which is the focus of this work. Our motivation is two-fold. On one hand, the existing work for sequential data procurement in the literature [?]do not work well for the situation of indoor localization; on the other hand, we want to build a concrete measurement of RSS data specifically for the active learning mechanism.

In this paper, we and propose a pricing mechanism according to the of the in order to get a higher . Our contributions are as following.

- We design an effective way to measure the what kind of RSS data should we purchase. In most cases, the collected RSS data is not idealy in the exact position, to design an effective way to measure the impact of those imperfect data is crucial to our system. We make a thorough analysis of the impact that imperfect data may exert on the result of localization through the probability model.
- We give the pricing strategy for the mechanism to acquire the high quality data. The mechanism has a theoretically better performance than the classical one proposed in **??**. The mechanism is robust in most indoor-localization situations, even the prior knowledge of the costs is not well understood and the noise in the crowdsensing data is rather arbitrary. And we further provide a most economical data procure mechanism when the target accuracy of RSS data is given.

The remaining of the paper is organized as following. The system structure and settings are given in section **??**. The measurement RSS data quality is presented in Section **??**. Section**??** gives the abstract definition for the online data procurement mechanism. Section **??** describes the mechanism under the simple assumption that the costs data are drawn in a distribution. Section **??** presents a generalized and more robust mechanism for the RSS data procurement. Section **??** gives our simulations and experiments for the mechanism we given before.

## II. RELATED WORK

### A. Fingerprinting based Indoor Localization

The RSS Fingerprinting based . The early technique,'nearest neighbour(s) in signal space'(NNSS) [?],outputs the location with minimum Euclidean distance between RSSes stored in the database and RSSes measured as the estimation of the users location. However, NNSS still dangles the possibility of accuracy enhancement in that it fails to realize the joint location estimation from multiple APs [?]. In order to take advantage of these multiple APs, Chintalapudi *et al.* bring up an algorithm entitled as EZ localization, whose main contribution occurs in estimating mobile devices without any pre-deployment support of multiple APs [?]. EZ will learn from those acquired fingerprints, which reflect the value of mean and standard deviation of the RSSes corresponding to different APs, during the collection phase. Wen *et. al* made a thorough study of the

### B. incentive mechanism for Indoor Localization

Fingerprinting based Indoor localization requires largePeng *et al.* bring up an incentive mechanism both stimulating data provision and ensuring high quality [?]. Jin *et al.* introduce a key metric, quality of information(QoI), which generally evinces the quality of users sensory data[?]. Taking QoI into consideration, the incentive mechanism can acquire data with higher quality making for further study like better identification for problems of medical devices [?].
In terms of new facets refreshing incentive mechanism research, Tham and Luo take timeliness of contributions into consideration [?].

### C. Online learning used for Crowdsourcing

## III. SYSTEM MODEL

In this section we describe our system model and give out the problem formulation.

We present a mobile crowdsensing system consisted of *RSS procurement mechanism*. For the purpose of performing accurate indoor localization in region $\mathcal{V}$, the data purchaser has to build the corresponding *Fingerprint Database* of *Received Signal Strength*(RSS). Therefore the data purchaser releases tasks of collecting data–RSS value on the platform. For a specific location $s \in \mathcal{V}$, we use $\mathbb{W}_s = \{w_1, w_2, ..., w_{N_s}\}$ to denote the corresponding applicants set. To simplify the notation, we omit the identification of $s$ in almost all the rest of this paper. Without loss of generality, we mainly focus on workers with the same location's data. It's rational that data purchaser need to buy several data points at one location since the RSS value is not constant, in fact it obeys some probability distribution, we assume that its probability density function is $\mathcal{D}(\cdot)$. Consequently, we need several amounts of samples to learn the distribution, more specifically, to estimate the mean value of RSS.

At the very beginning, the data purchaser needs to submit his *Pricing Mechanism* $\mathbb{M}$ to the platform. Here we consider the most nature trading scenario: these $N$ workers arrive in

a sequential way with his data $x_i$. Once agents $i$ arrives, he submit his bid $c_i$ to the platform and the platform compute its price $p_i$ using a mechanism $\mathbb{M}$. If $p_i \geq b_i$ then agents $w_i$ accept this transaction: the platform pays $b_i$ to worker $i$ and receives data $x_i$, otherwise the worker reject the transaction and the platform receives null signal.

#### TABLE I
#### NOTATIONS

| Notation | Remark |
|---|---|
| $\mathcal{V}$ | Indoor location region |
| $s$ | A specific location |
| $\mathbb{W}_s = \{w_1, \cdots, w_{N_s}\}$ | Applications set for location $s$ |
| $\mathcal{D}(\cdot)$ | Probability density function |
| $\mathbb{M}$ | Pricing mechanism |
| $w_i$ | The $i_{th}$ worker |
| $x_i, b_i, p_i$ | $w_i$'s data, bid, and corresponding price |

## IV. ANALYSIS OF TWO DIMENSION LOCALIZATION

### A. probability of error

The RSS value in the environment is hard to know, however, some research ?? has shown that the mean value and variance of the RSS follows a relative . Thus it is proper for us to make the assumption that the mean value of RSS in position $\vec{r}$ follows a continuous distribution, the Gaussian distribution,e.g. Since that some experiment result show that the RSS value may actualyFor a more general case, we may assume that the measured RSS value $P$ in location $\vec{r}$ follows the distribution of $f_{\vec{r}}(x; h)$, where the $h$ is the parameter of the distribution, or in other word, the hypothesis. As shown in Figure ??, we devide the physical space into many small circles each centered at $\vec{r}$ with radius $\vec{\delta}$, within each block we have a threshold $P_{high}$ and $P_{low}$ for the RSS value$P$. According to the $MLE$ principlr used in ??, which means that the probablity that the RSS falls into the ideal region must higher than the other reigon, $P_{high}$ and $P_{low}$ should satisfy that

$$
\begin{aligned}
f_{\vec{r}-\delta}(P_{high}) &= f_{\vec{r}}(P_{high}) \\
f_{\vec{r}+\delta}(P_{low}) &= f_{\vec{r}}(P_{low})
\end{aligned}
\tag{1}
$$

. Thus we may define the reliability as the probability of the system correctly estimate the user's location,

$$
R = \int_{P_{low}}^{P_{high}} f_r(P) dP
\tag{2}
$$

which means the probablity of RSS value tested in position $r$ lies within the interval $[P_{low}, P_{high}]$ However, in real circumstnaces, the $P_{high}$ and $P_{low}$ are acquired through the training process, during which may recieve the imperfect data and thus cause the result to be inaccurate.The speculated location of the RSS value may migrate from the original one. A very obvious situaiton is that the RSS value from one certain region may

falsely be recognized to from other area. We thus use this probability to define the error of the RSS data we collected.

$$P(error) = \iint_A f_A(Q)P(err|X = X_0)dxdy$$

$$= 2\iint_A f_A(Q)\int_{X_1}^{X_2}\frac{1}{\sqrt{2\pi}\sigma}e^{-\frac{(X-X_0)^2}{2\sigma^2}}dXdxdy \tag{3}$$

### B. The analysis of the loss

We can re write the probability of error as

$$P(error) = g(R, P) = h(|R - P|) \tag{4}$$

, where $R$ is the real RSS value and $\vec{P}$ is the measured data. In other word, $h(|R-P|)$ denotes the error between the real RSS value and the collected data. Obviously the real RSS value $x$ should satisfy that

$$R = arg\min_{x\in R} E_P h(|x - P|) \tag{5}$$

**Theorem 1.** *The expectation of the error $f(P; h)$ get its minimum when $P$ equals real RSS value $R$.*

Now we set $|R-P| = kln f(R, P)$ Assume that we sample $N$ data $x_1, ..., x_N$, we let $h*$ be the value that that minimize the $\frac{1}{N}\sum_{i=1}^N f(x_i; h)$. The theorem **??** above shows that when the number of data we collected is enough, the $t*$ we obtain from the data will approximate to the real RSS value $r$.

**Theorem 2.** *The average error of $\hat{t}$ obtained from collected sample has at least $1 - 2e^{-\frac{2\epsilon^2}{N}}$ the probablity that is within $\epsilon$ close to the average error of real RSS data $r$, that is*

$$Pr(\frac{1}{N}\sum_{i=1}^N f(x_i; \hat{t}) - E[f(x_i; r)] \leq \epsilon) \geq 1 - 2e^{-\frac{2\epsilon^2}{N}} \tag{6}$$

This theorem shows us that

### C. example of Gaussian Distribution

We may now give a more specific example of the theories we deduce above. We assume that the RSS value R folllow the Gaussian distribution with mean value of $\mu_r$, and according to many previous studies, we assume that the $\mu_{\vec{r}}$ is coninuous over $\vec{r}$. For the Gaussian distribution case, we give the specific form of these two threshold.

$$P_{high} = \frac{\mu(\vec{r}) + (\mu(\vec{r}) + \nabla\cos\varphi)}{2} = \mu(\vec{r}) + \frac{\nabla\cos\varphi}{2}$$
$$P_{low} = \frac{\mu(\vec{r}) + (\mu(\vec{r}) - \nabla\cos\varphi)}{2} = \mu(\vec{r}) - \frac{\nabla\cos\varphi}{2} \tag{7}$$

## V. LOW-COST DATA PURCHASING PROBLEM

In this section, we will give some preliminary about the task of our mechanism and prevailing principle used in statistical machine learning. We abstractly define the problem of the designing of the effective mechanism to acquire the RSS information collected by the crowds.

### A. Preliminaries

We first give the concept of loss and regret. The loss function that reflects the data quality is defined in the space $H \times Z \to R$, where $H$ is the hypothesis class and $Z$ is the space of the objects. We expect the loss function to get its minimum value when the data is exactly the ideal data. In our setting, the hypothesis $h$ is the mean value and variance of RSS fingerprinting data, and the hypothesis class $H$ is our expected internal of mean value and variance.

$$f_t(h_t) = \tag{8}$$

After we acquire the loss function, we give the concept of the regret function.

$$R(T) = \sum_{t=1}^T f_t(h_t) - \min_{h^*\in H}\sum_{t=1}^T f_t(h_t^*) \tag{9}$$

where $h^*$ is the optimal choice, causing the least loss in our solution space $H$. The regret function reflexes how the data deviate from the desired value, the real mean and variance of RSS. We also make some assumptions for this problem

1) the agents have nothing to do with the costs
2) the
3) the

### B. Online Gradient Descent

Online learning is a widely used learning paradign. The goal of online learning is to produce the best hypothesis when data is in sequential order We here use the classical Online Gradient Descent(OGD) algorithm to work as the Online Algorithms. It has been proved that the OGD has an upper bound of regret of $O(\sqrt{T})$, which ensures that the average regret tends to zero when $T$ goes to infinite. There are also many kinds of other Online Algorithms which can be found in refOLServey, etc. The OGD is described as following. In each time $t$, we obtain a $h_t$ according to

$$h_t = h_{t-1} - \eta\nabla f_t(h_{t-1}) \tag{10}$$

.

### C. Importance Weighting technique

In tradational online learning problem, all the data will be used to produce the total regret. In our low-cost purchasing problem, the mechanism do not get access to data and obtain a loss in each time $t$. the estimation of loss is $E(\sum_{t=0}^T \delta_t f_t) = \sum_{t=0}^T q_t f_t$, where $\delta_t$ is the function showing whether the data is procured. However, the definition of regret in [**??**] still includes all the loss in each time $t$, whether it has been used or not. In order to get an unbiased estimator of the regeret, we define

$$\hat{f}_t(h) = \begin{cases} \frac{f_t(h_t)}{q_t} & data \quad access \quad to \quad RPM \\ 0 & else \end{cases} \tag{11}$$

.

### D. Problem definition

We consider that the data collected through crowdsensing coming in a sequence of $d_1, ,, d_T$, with each of them contains a cost $c_1, ,, c_T$. We should design a pricing mechanism that can choose how much we should pay for the data. However, we have no means to know either the quality of data is good enough for localization or there will be a better one coming in the sequence. Under the framework of online machine learning, we formally define our RSS data Procure Mechanism ($RPM$) is defined as following.

**Definition 1.** *Given a sequence of data $d_1, ..., d_T$ coming in time $1, ,,,, T$ with each data possessing a posted price $c_t$, $c_t \in [0, M]$.*

1) *The RPM post a hyphhesis $h_t$ from OGD*
2) *The RPM post a price $p_t$ according to a distribution $G_t$ over $[0, M]$.*
3) *If the $p_t > c_t$ agent accepted the price, the RPM send the loss function $f(h_t)/q_t$ back to the OGD and pay for the posted price $p_t$. If $p_t < c_t$ the agent rejected the price, the mechanism send a null data to the OGD .*

*The mechanism outputs a final hypothesis $\overline{h} \in H$*

And clearly, our kernel problem is to find the best distribution $G_t$ used for the mechanism to post its price

### E. online batch to conversion

The meachanism and online learning algorithm produces a sequence of hypothesis $h_1, ,, h_T$. The main goal of our algorithm is to get the best hypothesis $\overline{h}$, the mean value and variance of RSS, from the sequence. One simple aproach is to average every hypothesis $h_t$ acquired in each time $t$.

$$\overline{h} = \sum_{t=1}^{n} h_t \tag{12}$$

It has been proved that ...,

## VI. THE MAIN SETTING: REGRET MINIMIZATION SENARIO

In this senario, the mechanism has a fixed budget. The target of the mechanism is that in each round $t$, the produce the minimum regret defined in **??** . In this section, we will give the exact form of the distribution $G_t$ and the analysis of the regret bound according to this distribution.

### A. Estimate the upper bound of regret

We will find the upper bound of the regret defined in (**??**) produced by RPM. In normal case, the regret bound of OGD is

$$\frac{||h||^2}{2\eta} + \eta \sum_{t=1}^{T} \nabla f_t(h_t)^2 \tag{13}$$

, which is a well known result. Under the importance weighting framework, we give the regret bound in the following lemma

**Lemma 1.** *The regret bound produced by RPM in* **??** *is bounded by*

$$R(h) \leq \frac{||h||^2}{2\eta} + \eta E(\sum_{t=1}^{T} \frac{\nabla f_t(h_t)^2}{q_t}) \tag{14}$$

. The **??** is quite easy to be proved under our setting that the loss function $f_t$ is of strong convexity.

### B. Derivation of the Regeret Minimization Problem

In each time $t$, the RPM need to post a price $p$ according to a distribution $g$ in order to get a minimum regret, we thus reduce the problem of designing a mechanism into an optimization problem

$$\min \sum_{i=1}^{n} \frac{\nabla f_i^2}{1 - F_i(c_i)}$$
$$s.t. \quad \sum_{i=1}^{n} \int_{c_i}^{M} x dF_i(x) \leq B \tag{15}$$

where $\forall c_i, 0 \leq c_i \leq M$, and $F(0) = 0, F(M) = 1$

**Theorem 3.** *The optimal solution of the optimization problem [**??**] is in the form*

$$F_t(c) = \begin{cases} 1 - \frac{\nabla f_t}{\sqrt{\lambda c - \beta}} & c \in (\frac{\nabla f_t^2 + \beta}{\lambda}, M] \\ 0 & else \end{cases} \tag{16}$$

*Proof:* We first give our function space $V = \{y | y(0) = 0, y(M) = 1\}$. And we denote our cost function as

$$M(F_1, ,, F_n) = \sum_{i=1}^{n} \frac{\alpha_i}{1 - F_i(c_i)}.$$

Then the augmented Lagrange function is derived as

$$J(F_1, ,, F_n, \lambda) = M(F_1, ,, F_n) + \lambda(\sum_{i=1}^{n} \int_{c_i}^{M} x dF_i(x) - B)$$

According to the Gateaux Deravative, we obtain that for $\forall \hat{F} \in V$

$$\delta J|_{F_t}(\hat{F}_t - F_t) = \int_{c_t}^{M} (-\frac{\alpha_t}{(1 - F_t(c_t))^2} + \lambda x)(\hat{f}(x) - f(x)) dx$$

if $\overline{F}$ is the local minimum, then we have

$$\delta J(\hat{F}_t - \overline{F_t}) \geq 0$$

holds for every $\hat{F} \in V$. Noticing that

$$\int_{0}^{M} f_t(x) - f(x) dx = 0$$

We must have

$$-\frac{\alpha_t}{(1 - F_t(c_t))^2} + \lambda x \geq 0$$

hold on every where on $[c_t, M]$ thus we obtain that

$$F_t(c) = \begin{cases} 1 - \frac{\nabla f_t}{\sqrt{\lambda c - \beta}} & c \in (\frac{\nabla f_t^2 + \beta}{\lambda}, M] \\ 0 & else \end{cases} \tag{17}$$

One major difficulty for the is the determination of the determination of $\beta$ and $\lambda$ Noticing that $F(x)$ is not continuous, according to Stieltjes Integral, we rewrite the constraint as following

$$
\begin{aligned}
&\sum_{t=1}^{T}(\int_{c_t}^{M} x dF_t(x)) \\
=&\sum_{t=1}^{T}(\int_{c_t}^{M} x f_i(x)dx + (1-F_i(M)M) \\
\leq&\sum_{t=1}^{T} \nabla f_t(\frac{2}{\lambda}\sqrt{\lambda M - \beta} + \frac{c_t}{\sqrt{\lambda c_t - \beta}} - \frac{2}{\lambda}\sqrt{\lambda c_t - \beta}) \\
\leq& B
\end{aligned}
$$

The Stieltjes Integral here has its practical significance. Because we assume that the cost lies between $[0, M]$, in other word, the mechanism do not accept any price higher than $M$, thus for all posted price $c$ that are higher than $M$, the mechanism will only pay $M$ instead of $c$.

Now since we get the solution of the $F_t$, the remaining work is to determine the parameters $\lambda$ and $\beta$, we go back to our initial optimizationproblem that minimize the regret bound. The Lagrangian is thus given as follows

$$
\begin{aligned}
L(\mu, \beta, \lambda) = \sum_{t} \Big( &\nabla f_t\Big(\sqrt{\lambda c_t - \beta} + \mu(\frac{2}{\lambda}\sqrt{\lambda M - \beta} \\
&+ \frac{c_t}{\sqrt{\lambda c_t - \beta}} - \frac{2}{\lambda}\sqrt{\lambda c_t - \beta})\Big)\Big) - \mu B
\end{aligned} \tag{18}
$$

According to the complementary relaxation condition, $\mu \neq 0$, which means that the constraint condition in **??** for the optimal solution is strict. To get the analytic solution of the optimal value of $\beta$ and $\lambda$ is infeasible , thus we use the numeric solution for the equation **??**.

*C. Analysis of the result*

We first simply set the $\beta = 0$ for a special case. Through simple calculation, we can have an estimation of $\lambda_0$ as following

$$
\lambda_0 = \frac{T}{B}(2\theta_0 - \theta) \tag{19}
$$

where $\theta_0 = \frac{1}{T}\sum_t \nabla f_t \sqrt{M}$, $\theta = \frac{1}{T}\sum_t \nabla f_t \sqrt{c_t}$ Since that $\partial L/\partial \beta > 0$, $\partial L/\partial \lambda < 0$, we obtain that the optimal solution $(\beta^*, \lambda^*)$ statisfy that $\beta^* > \beta^0$, $\lambda^* < \lambda_0$ Thus we have the estimate of the upper bound of the regret of RPM in theorem **??**.

**Theorem 4.** *The regret of RPM in produced by the algorithm in **??** is bounded by*

$$
Regret < O(\frac{T}{\sqrt{B}}(2\theta_0 - \theta)\alpha) \tag{20}
$$

*where $\alpha = O(\sqrt{1 - \frac{\beta B^2}{T\theta^2}})$*

*Proof:* To be determined ∎

One problem in this situation is that we may not get enough prior knowledge to both $c_t$ and $\nabla f_t$. One way to solve the problem is that we initially set $\beta$ to a fixed value and $\lambda$ to a very small value, e.g. 0.0001. Then in each time $t$, we update the value of $\lambda$ with

$$
\theta_0^{(t)} = \sum_{i=1}^{t-1} \frac{\nabla f_t(h_t)}{t-1}\sqrt{M} \tag{21}
$$

$$
\lambda^{(t)} = \frac{T^2}{B^2 M}\theta_0^2 + \frac{\beta}{M} \tag{22}
$$

## VII. THE BUDGET MINIMIZATION SENARIO

In this senario, we consider the problem to find the most money-saving way to acquire the data inorder to achieve a satisfactory regret bound $R$. The objective function is in the form of a integral, which is not an easy problem of the classical optimization problem. Besides, to solve the very exact form of the budget . Thus we use the of the budget B as following

$$
\sum_t c_t q_t \leq B \leq \sum_t M q_t \tag{23}
$$

Thus what we have to do is to solve the optimization problem of the form

$$
\begin{aligned}
&\min_{q_t} c_t \\
&s.t. \sum_t \frac{\nabla f_t^2}{q_t} \leq R \\
&0 \leq q_t \leq 1
\end{aligned} \tag{24}
$$

*A. The optimal mechanism*

Consider the convexity of the objective function, we give the Lagrangian

$$
L = \sum_t c_t q_t - \lambda(-\sum_t \frac{\nabla f_t^2}{q_t} + R - \sum_t \mu_t(1-q_t)) \tag{25}
$$

The optimal K-T condition of the problem **??** is

$$
\frac{\partial L}{\partial q_t} = c_t - \lambda[\frac{\nabla f_t^2}{q_t^2}] - \mu_i = 0 \tag{26}
$$

when $q_t = 1$, we get $u_i \neq 0$, when $q_t \neq 1$, $\mu_i = 0$, thus we have

$$
q_t = \min\{1, \sqrt{\frac{\lambda}{c_t}}\nabla f_t\} \tag{27}
$$

According to our constraint condition

$$
\sum_t \sqrt{\frac{c_t}{\lambda}}\nabla f_t \leq R \tag{28}
$$

we can get an approximation of the $\sqrt{\lambda}$ through simple calculation

$$
\sqrt{\lambda} = \frac{T}{R}\theta \tag{29}
$$

where we use $\theta$ to denote the term $\frac{1}{T}\sum_t \sqrt{c_t}\nabla f_t$. Since **??** holds for $\forall c_t$, and $c_t$ is arbitrarily given.We may assume that the convoluted distribution function of the price mechanism is of the form

$$
F_t(c) = 1 - \sqrt{\frac{\lambda}{c}}\nabla f_t \tag{30}
$$

And the PDF is

$$f(c) = \frac{1}{2}\sqrt{\frac{\lambda}{c^3}}\nabla f_t \tag{31}$$

### B. result analysis

Now we can make a relatively more precise estimate the budget $B$

$$E(B) = \sum_t \int_{c_t}^{M} cf(c)dc \tag{32}$$

$$= \frac{T^2}{R}\theta\varphi \tag{33}$$

where $\varphi = \sum_t \frac{1}{T}\nabla f_t(\sqrt{M} - \sqrt{c_t})$

## VIII. Experiments and Simulations

In this section, we conduct the experiments and simulations to validate the ,,,of our system model and data procurement mechanism. We use the data collected in FoxCom Shanghai,where we tested the RSS value of 10 AP in 13 different locations, the distribution of the location and AP points are shown in Figure**??**. Since that the workers collected our data did not ased for reward to us, we simply simulate the costs of the data through a normal distribution with mean value of $0.5$ and variance of $1$.
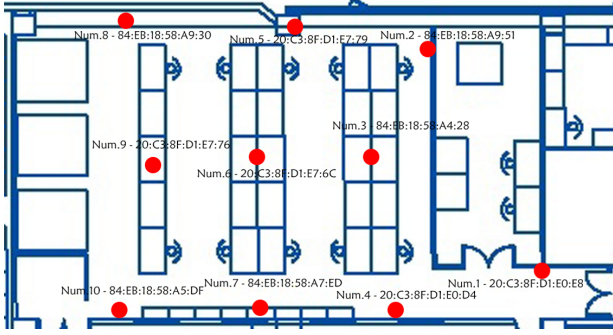


Fig. 1. Location Description.

## IX. Conclusion and Future Work

### References

[1] Z. Yang, Z. Zhou and Y. Liu, "From RSSI to CSI: Indoor localization via channel response," *ACM Comput. Surv.*, vol. 46, no. 2, pp.1-32, 2013.
[2] K.Kaemarungsi, and P.Krishnamurthy, "Modeling of indoor positioning systems based on location fingerprinting," in *Twenty-third AnnualJoint Conference of the IEEE Computer and Communications Societies*, 2004, vol. 2, pp. 1012–1022.
[3] Y. Wen, X. Tian, X. Wang and S. Lu, "Fundamental limits of RSS fingerprinting based indoor localization," in *Proc. IEEE INFOCOM*, 2015, pp.2479–2487.
[4] R.K.Ganti, F.Ye, H,Lei, "Mobile crowdsensing: current state and future challenges," *IEEE Communications Magazine*, vol.49, no.11, pp.32-39, 2011.
[5] S.Hu, L.Su, H.Liu, H.Wang, and T.F.Abdelzaher, "Smartroad: Smartphone-based crowd sensing for traffic regulator detection and identification," *ACM Transactions on Sensor Networks*, vol.11, no.4, pp.55, 2015.
[6] M.Mun, S.Reddy, K.Shilton, N.Yau, J.Burke, D.Estrin, M.Hansen, E.Howard, R.West, and P.Boda, "PEIR, the personal environmental impact report, as a platform for participatory sensing systems research," in *Proc.ACM MobiSys*, 2009, pp.55–68.
[7] R.Rana, C.Chou, S.Kanhere, N.Bulusu, and W.Hu, "Earphone:An end-to-end participatory urban noise mapping," in *Proc.ACM/IEEE IPSN*, 2010, pp.105–116.
[8] R.Gao, M.Zhao, T.Ye, F.Ye, Y.Wang, K.Bian, T.Wang, and X.Li, "Jigsaw: Indoor floor pan plan reconstruction via mobile crowdsensing," in *Proc.ACM MobiCom*, 2014, pp.249–260.
[9] Y.Wen, J.Shi, Q.Zhang, X.Tian, Z.Huang, H.Yu, Y.Cheng, and X.Shen, "Quality-driven auction-based incentive mechanism for mobile crowd sensing," *IEEE Transactions on Vehicular Technology*, vol.64, no.9, pp.4203-4214, 2015.
[10] M.Zinkevich, "Online convex programming and generalized infinitesimal gradient ascent," School of Computer Science, Carnegie Mellon University, 2003.
[11] S.Shalev-Shwartz, "Online learning and online convex optimization," *Foundations and Trends in Machine Learning*, vol.4, no.2, pp.107-194, 2011.
[12] J.Abernethy, Y.Chen, C.J.Ho, and B.Waggoner, "Low-cost learning via active data procurement," in *Proceedings of the Sixteenth ACM Conference on Economics and Computation*, 2015, pp.619–636.
[13] A.Beygelzimer, S.Dasgupta, and J.Langford, "Importance weighted active learning," in *Proceedings of the 26th Annual International Conference on Machine Learning*, 2009, pp.49–56.
[14] D.Liberzon, "Calculus of variations and optimal control theory: a concise introduction," Princeton University Press, 2012.
[15] A.Roth, and G.Schoenebeck, "Conducting truthful surveys, cheaply," in *Proceedings of the 13th ACM Conference on Electronic Commerce*, 2012, pp.826–843.