

# A Clonal Selection Algorithm with Levenshtein Distance based Image Similarity in Multidimensional Subjective Tourist Information and Discovery of Cryptic Spots by Interactive GHSOM

Takumi Ichimura

Faculty of Management and Information Systems,  
Prefectural University of Hiroshima  
1-1-71, Ujina-Higashi, Minami-ku,  
Hiroshima, 734-8559, Japan  
Email: ichimura@pu-hiroshima.ac.jp

Shin Kamada

Graduate School of Comprehensive Scientific Research,  
Prefectural University of Hiroshima  
1-1-71, Ujina-Higashi, Minami-ku,  
Hiroshima, 734-8559, Japan  
Email: shinkamada46@gmail.com

**Abstract**—Mobile Phone based Participatory Sensing (MPPS) system involves a community of users sending personal information and participating in autonomous sensing through their mobile phones. Sensed data can also be obtained from external sensing devices that can communicate wirelessly to the phone. Our developed tourist subjective data collection system with Android smartphone can determine the filtering rules to provide the important information of sightseeing spot. The rules are automatically generated by Interactive Growing Hierarchical SOM. However, the filtering rules related to photograph were not generated, because the extraction of the specified characteristics from images cannot be realized. We propose the effective method of the Levenshtein distance to deduce the spatial proximity of image viewpoints and thus determine the specified pattern in which images should be processed. To verify the proposed method, some experiments to classify the subjective data with images are executed by Interactive GHSOM and Clonal Selection Algorithm with Immunological Memory Cells in this paper.

**Index Terms**—Levenshtein Distance, Clonal Selection Algorithm, Image Analysis, Immunological Memory Cells, Growing Hierarchical SOM, Interactive GHSOM, Smartphone based Participatory Sensing System, Tourist Informatics, Knowledge Discovery

## I. INTRODUCTION

The current information technology can collect various data sets because the recent tremendous technical advances in processing power, storage capacity and network connected cloud computing. The sample record in such data set includes not only numerical values but also language, evaluation, and binary data such as pictures. The technical method to discover knowledge in such databases is known to be a field of data mining and developed in various research fields.

Mobile Phone based Participatory Sensing (MPPS) system involves a community of users sending personal information

and participating in autonomous sensing through their mobile phones [1]. Sensed data can be obtained from sensing devices present on mobiles such as audio, video, and motion sensors, the latter available in high-end mobile phones. Sensed data can also be obtained from external sensing devices that can communicate wirelessly to the phone. Participation of mobile phone users in sensorial data collection both from the individual and from the surrounding environment presents a wide range of opportunities for truly pervasive applications. The tourist subjective data collection system with Android smartphone has been developed [2]. The application can collect subjective data such as pictures with GPS, geographic location name, the evaluation, and comments in real sightseeing spots where a tourist visits and more than 500 subjective data are stored in the database. Attractive knowledge discovery for sight seeing spots is required to promote the sightseeing industries.

We have already proposed the classification method from the collected subjective data by the interactive GHSOM [3], [4] and the knowledge is extracted from the classification results of the interactive GHSOM by C4.5 [5]. However, the image data was not included in the classification tasks, because it is too large amount of information to realize the extraction of specified characteristics from images.

There is currently an abundance of vision algorithms which are capable of determining the relative positions of the viewpoints from which the images have been acquired. However, very few of these algorithms can cope with unordered image sets for which no a prior proximity ordering information is available. Image localization can be addressed in the framework of the fundamental structure and motion (SaM) estimation problem and benefits from the wide field of view offered by Smartphone camera. This is because a wide field of view facilitates capturing large portions of the environment with few images and without resorting to the use of movable gaze control mechanisms such as pan-tilt units. Furthermore, environment features remain visible in large subsets of images

©2013 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

and critical surfaces are less likely to cover the whole visual field.

The definition of relative positions and orientations of the viewpoints corresponding to a set of unordered central images is an important procedure to be statistical analysis in image retrieval. The idea in the proposed approach employs the Levenshtein distance[6] to deduce the spatial proximity of image viewpoints and thus determine the specified pattern in which images should be processed. Horizontal matching method for localizing unordered panoramic images has been proposed[7]. In the method, all images have been acquired from a constant height above a planar ground and operates sequentially by the Levenshtein distance. Our proposed method can process not only in horizontal matching but also in vertical matching. In this paper, the photographs are divided into some categories according to the similarity by the clonal selection algorithm with immunological memory cells before the classification by GHSOM.

The area of artificial immune system (AIS) has been an ever-increasing interested in not only theoretical works but applications in pattern recognition, network security, and optimization [8], [9]. AIS uses ideas gleaned from immunology in order to develop adaptive systems capable of performing a wide range of tasks in various research areas. Gao indicated the complementary roles of somatic hypermutation (HM) and receptor editing (RE) and presented a novel clonal selection algorithm called RECSA model by incorporating the Receptor Editing method [10]. The immunological memory which leads to a perception that an individual is immune to a particular agent is realized by the clustering of the generated antibodies[11].

The remainder of this paper is organized as follows. In Section II, the clonal selection theory with memory cells will be explained briefly. The idea about the antibody structure of images by Levenshtein Distance and experimental results are discussed in Section III. Section IV describes the algorithm of interactive GHSOM and its interface tool. Section V explains the tourist subjective data and the experimental results. In Section VI, we give some discussions to conclude this paper.

## II. CLONAL SELECTION ALGORITHM WITH IMMUNOLOGICAL MEMORY

Clonal Selection Algorithm with Immunological Memory(CSAIM) model has been proposed to introduce an idea of immunological memory into the RECSA model. This section describes the structure of antibody in RECSA model to the medical diagnosis briefly. The further details about the CSAIM algorithm was described in [11].

### A. Antibody for Classification Problem

This subsection describes the antibody for classification problem about the structure, the method of somatic hypermutation and receptor editing, and affinity.

#### 1) Structure of Antibody for Classification Problem:

Fig.1 shows the structure of antibody in the classification problem[11].  $w_k$ ,  $\theta$  is the weight of antibody and threshold,

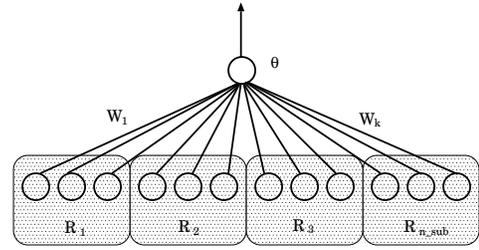


Fig. 1. The antibody structure

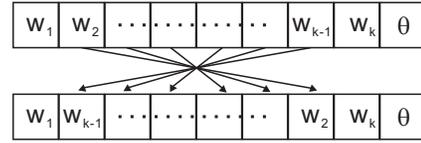


Fig. 2. RE for  $w_2, w_{k-1}$

respectively.  $R_1, \dots, R_{n_{sub}}$  indicate the sub-region in the problem, because some classification problem can be divided into  $n_{sub}$  sub tasks. That is, a region is expert for the specified task in classification.

2) *Somatic Hypermutation and Receptor Editing*: HM updates the randomly selected  $w_i$  and  $\theta$  for a paratope  $P = (w_1, \dots, w_k, \theta)$  as follows.

$$w_i = w_i + \Delta w, \theta = \theta + \Delta \theta,$$

where  $\Delta w$ ,  $\Delta \theta$  are  $-\gamma_w < \Delta w < \gamma_w$ ,  $-1 < \Delta \theta < \gamma_\theta$ , respectively.  $\gamma_w$  and  $\gamma_{theta}$  are a small number.

RE makes a crossover of 2 set of  $w_i$  for a paratope as shown in Fig.2.

3) *Affinity*: The system calculates the degree of affinities between antibody and antigen by using Eq.(1) and Eq.(2).

$$f(x^p) = \begin{cases} 1 & \text{if } |\sum_{i=1}^k w_i x_i^p - \theta| \geq E_{sim} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

$$g(x^p) = \begin{cases} 1 & \text{if } f(x^p) = x_{Target}^p \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Eq.(3) calculates the degree of affinity.

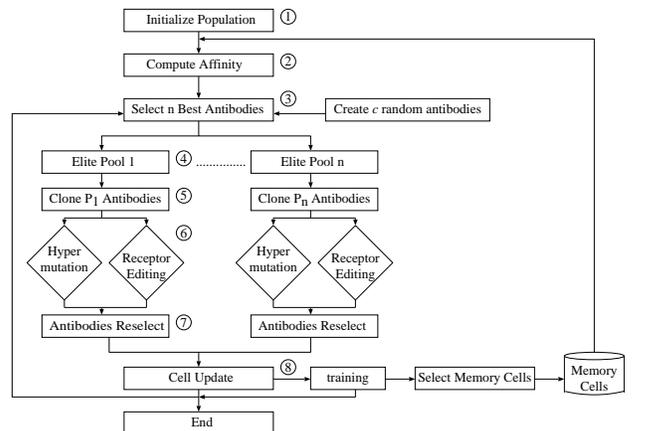


Fig. 3. A flow of CSAIM model

$$Affinity = \sum_{p=1}^{tr\_num} g(x^p), \quad (3)$$

where  $x^p$  means the  $p$ th sample in  $tr\_num$  training cases and  $x_{Target}^p$  is 1 if an example is a Target, otherwise 0.

### B. RECSA Model[10]

The shape-space model aims at quantitatively describing the interactions among Ags and Abs (Ag-Ab) [12]. The set of features that characterize a molecule is called its generalized shape. The Ag-Ab codification determines their spatial representation and a distance measure is used to calculate the degree of interaction between these molecules.

The Gao's model [10] can be described as follows.

- 1) Create an initial pool of  $m$  antibodies as candidate solutions ( $Ab_1, Ab_2, \dots, Ab_m$ ).
- 2) Compute the affinity of all antibodies: ( $D(Ab_1), D(Ab_2), \dots, D(Ab_m)$ ).  $D()$  means the function to compute the affinity.
- 3) Select  $n$  best individuals based on their affinities from the  $m$  original antibodies. These antibodies will be referred to as the elites.
- 4) Sort the  $n$  selected elites in  $n$  separate and distinct pools in ascending order. They will be referred to as the elite pools.
- 5) Clone the elites in the pool with a rate proportional to its fitness. The amount of clone generated for these antibodies is given by Eq.(4).

$$P_i = \text{round}\left(\frac{n-i}{n} \times Q\right), \quad (4)$$

where  $i$  is the ordinal number of the elite pools,  $Q$  is a multiplying factor for determining the scope of the clone and  $\text{round}()$  is the operator that rounds towards the closest integer. Then, we can obtain  $\sum P_i$  antibodies as  $((Ab_{1,1}, Ab_{1,2}, \dots, Ab_{1,p_1}), \dots, (Ab_{n,1}, Ab_{n,2}, \dots, Ab_{n,p_n}))$ .

- 6) Subject the clones in each pool through either hypermutation or receptor editing process. The mutation rates,  $P_{hm}$  for hypermutation and  $P_{re}$  for receptor editing given by Eq.(5) and Eq.(6), are inversely proportional to the fitness of the parent antibody,

$$P_{hm} = a/D() \quad (5)$$

$$P_{re} = (D() - a)/D(), \quad (6)$$

where  $D()$  is the affinity of the current parent antibody and  $a$  is an appropriate numerical value.

- 7) Determine the fittest individual  $B_i$  in each elite pool from amongst its mutated clones. The  $B_i$  is satisfied with the following equation.

$$\begin{aligned} D(B_i) &= \max(D(Ab_{i,1}), \dots, D(Ab_{i,p_i})), \\ i &= 1, 2, \dots, n \end{aligned} \quad (7)$$

- 8) Update the parent antibodies in each elite pool with the fittest individual of the clones and the probability  $P(Ab_i \rightarrow B_i)$  is according to the roles: if  $D(Ab_i) < D(B_i)$  then  $P = 1$ , if  $D(Ab_i) \geq D(B_i)$  then  $P = 0$ , if  $D(Ab_i) \geq D(B_i)$ ,  $i \neq 1$  then  $P = \exp\left(\frac{D(B_i) - D(Ab_i)}{\alpha}\right)$ .

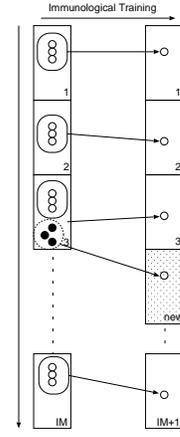


Fig. 4. A clustering method of memory cells

- 9) Replace the worst  $c (= \beta \times n, \beta$  is the parameter.) elite pools with new random antibodies once every  $t$  generations to introduce diversity and prevent the search from being trapped in local optima.
- 10) Determine if the maximum number of generation  $G_{max}$  to evolve is reached. If it is satisfied with this condition, it terminates and returns the best antibody. Otherwise, go to Step 4).

### C. Immunological Memory Cell

Clustering Memory Cells are required to classify the antibodies responding the specified samples. This paper realizes the clustering by allocating the generated antibodies by RECSA model into some categories. The initial number of categories is predefined and a new category is created according to training situation. Fig. 4 shows the clustering method of memory cells. Similar antibodies crowd around an appropriate point in each category, and then only central antibody of the crowd can become a memory cell. However, we may meet that memory cells can not recognize some of samples in the data set. In such a case, some new generated antibodies by RECSA model tries to respond to the mis-classification of the samples, if the similar antibodies make a crowd.

To find the crowd of similar cases, the system checks whether the Euclidean distance between normalized training sample and its corresponding antibody is smaller than the predetermined parameter  $\mu_\theta$ .

The similarity is measured by the following. Let  $\vec{d} = (d_1, \dots, d_i, \dots, d_k)$  be the elements of input signal and  $\vec{h} = (h_1, \dots, h_i, \dots, h_k)$  be the element of antibody.

In order to calculate the distance between the sample and the antibody, the range of sample is changed to that of antibody as follows.

$$d'_i = d_i \times \frac{h_j}{d_j} (d_i \neq 0 \wedge h_i \neq 0),$$

where  $d_j$  is the min value of element in the input sample.

Then, if the Euclidean distance between  $\vec{d}'$  and  $\vec{h}$  is smaller than  $\mu_\theta$ , the antibody can respond the sample. In this paper,  $\mu_\theta$  is the summation of 12 input elements.



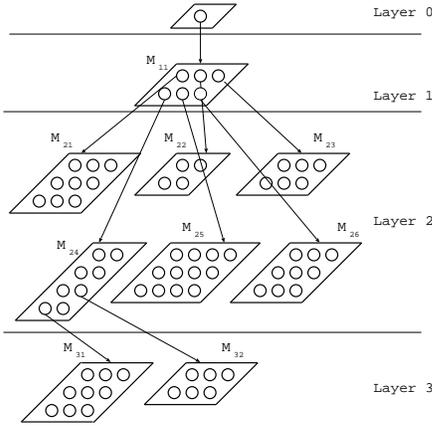


Fig. 7. A Hierarchy Structure in GHSOM

We proposed the reconstruction method of hierarchy of GHSOM even if the deeper GHSOM is performed[3], [4]. A stopping criterion for stratification is defined. Moreover, if the quantization error is large and the condition of hierarchies is not satisfied, the requirements for redistribution of error are defined.

Case1) If Eq.(9) and the classification capability for hierarchies are satisfied, stop the process of hierarchies and insert new units in the map again.

$$n_k \leq \alpha n_I, \quad (9)$$

where,  $n_k$ ,  $n_I$  mean the number of input samples for the winner unit  $k$  and of the all input samples  $I$ , respectively. The  $\alpha$  is a constant.

Case2) If the quantization error is not larger and the addition of layer is not executed, we may meet the situation that the quantization error of a unit is larger than the quantization error in an overall map. If Eq.(10) is satisfied, then a new unit is inserted.

$$qe_k \geq \beta \tau_1 \sum_{y \in \mathbf{S}_k} qe_y, \quad (10)$$

where  $\mathbf{S}_k$  is the set of winner units  $k$ .  $\beta \tau_1$  is a constant for the quantization error.

### C. An interface of interactive GHSOM

We developed the Android smartphone based interface of interactive GHSOM to acquire the knowledge intuitively. This tool was developed by Java language. Fig. 8 shows the clustering results of Iris data set [16] by GHSOM. The notation  $[R][01][10] : 11$  as shown in Fig. 8 represents the location of unit in the connection from the top level  $[R]$ .  $[R]$  means a root node. The numerical value('11') shows the number of samples divided into the leaf map after the sequence of classification  $[R][01][10]$ . The numerical values in the brackets mean the position of units in the corresponded map. The first letter (e.g. '0') and the second letter (e.g. '1') are the position in the column and the row in the map(e.g. '01'), respectively.

The similar color of units represents an intuitive understanding of similar pattern of samples. If the number of units in a map are increased, only a few samples could be classified into a new generated unit. Once the unit connected to the map is

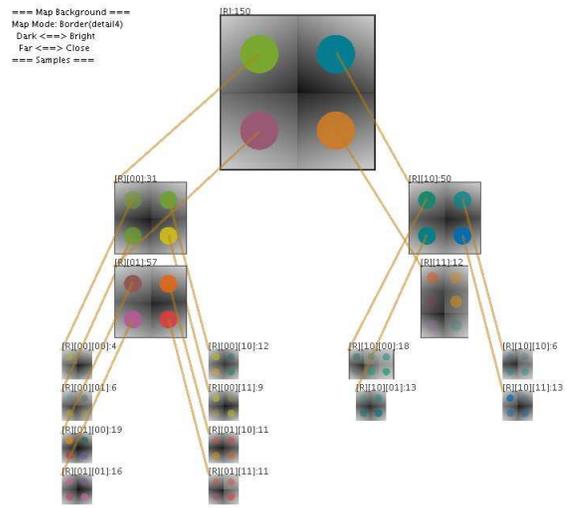


Fig. 8. Simulation Result for Iris data set

selected, the method re-calculates to find an optimal set of weights in the local tree structure search and then a better structure is depicted.

When the corresponding unit is touched, the system calculates the 4 allocated samples. The system continues to classify again till the user determines the GHSOM structure. We call the method the interactive GHSOM in the interactive process. The calculation result by the interactive GHSOM as shown in Fig. 8 is obtained and the effectiveness of the interactive GHSOM is shown as results of empirical studies.

The classification by GHSOM shows the tree structure of clusters and the connection among them. The detailed knowledge cannot be represented in the form of If-Then rules. Despite low resolution in knowledge representation, we grasp the rough sketch of knowledge structure because GHSOM shows the samples divided into each unit on the map as shown in Fig. 8. Moreover, there is only a few samples in each unit. Therefore, knowledge discovery is executed by grasping the structure and a grain of knowledge.

## V. EXPERIMENTAL RESULT

Participation of mobile phone users in sensorial data collection both from the individual and from the surrounding environment presents a wide range of opportunities for truly pervasive applications [1]. Our developed Android smartphone application [2] can collect the tourist subjective data in the research field of MPPS. The collected subjective data consist of jpeg files with GPS, geographic location name, the evaluation of  $\{0, 1, 2, 3, 4\}$  and comments written in natural language at sightseeing spots to which a user really visits. The 'comment' must be converted the number of words extracted from html files in the Tourist websites to a numerical value. The term frequency in the subjective comments is calculated by TF-IDF (term frequency inverse document frequency) method [17]. More than 500 subjective data are stored in the database through MPPS.

## VI. CONCLUSIVE DISCUSSION

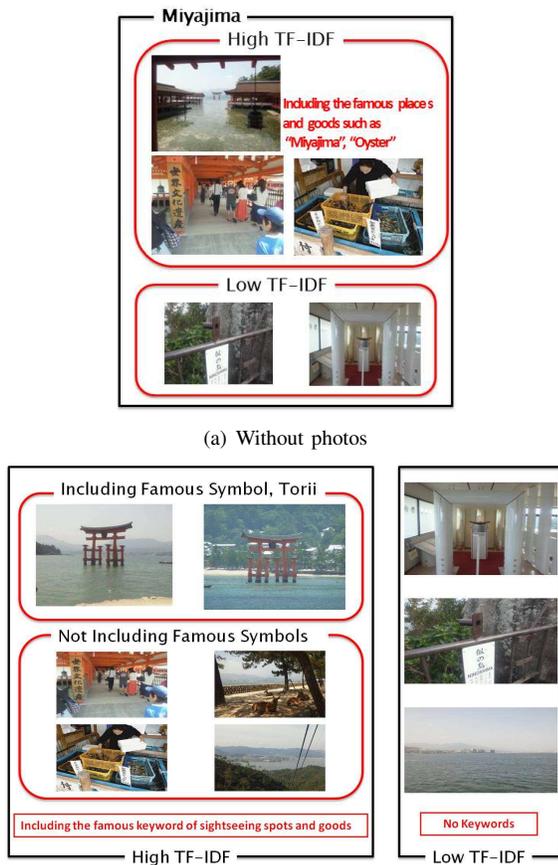
This paper presents an efficient and robust method for extracting the specified images and arrange them in ascending order of the similarity, which represents the degree including the image of landmark. Matching a limited amount of image data has been shown to suffice for registering the images. CSAIM method can classify the specified pattern extracted from the images and then the splices in the pattern are divided into some groups. The tourist subjective data with the specified image which are collected through MPPS is classified by Interactive GHSOM. As a result, the discovery of cryptic sightseeing spots is executed. The similarity to the unfamiliar landmark in sightseeing spots does not measured. In order to improve such a problem, more subjective data is required to classify them and the social action among participants will be investigated in future.

### ACKNOWLEDGMENT

This work was supported by JSPS KAKENHI Grant Number 25330366.

### REFERENCES

- [1] N.D.Lane, E.Miluzzo, L.Hong, D.Peebles, T.Choudhury, A.T.Campbell, *A survey of mobile phone sensing*, IEEE Communications Magazine, Vol.48, No.9, pp.140-150, 2010.
- [2] ITProducts, *Hiroshima Tourist Map*, <https://market.android.com/details?id=jp.itproducts.KankouMap>, [online], 2011.
- [3] T.Ichimura, S.Kamada, and K.Kato, *Knowledge Discovery of Tourist Subjective Data in Smartphone Based Participatory Sensing System by Interactive Growing Hierarchical SOM and C4.5*, Intl. J. Knowledge and Web Intelligence, Vol.3, No.2, pp.110-129, 2012.
- [4] T.Ichimura, S.Kamada, *A Generation Method of Filtering Rules of Twitter Via Smartphone Based Participatory Sensing System for Tourist by Interactive GHSOM and C4.5*, 2012 IEEE International Conference on Systems, Man, and Cybernetics (IEEE SMC 2012), pp.110-115, 2012.
- [5] J.R.Quinlan, *Improved use of continuous attributes in c4.5*, Journal of Artificial Intelligence Research, No.4, pp.77-90, 1996.
- [6] V.I. Levenshtein, *Binary codes capable of correcting deletions, insertions and reversals*, Soviet Physics - Doklady, Vol.10, No.8, pp.707710, 1996.
- [7] D.Michel, A.A.Argyros, M.I.A.Lourakis, *Horizon matching for localizing unordered panoramic images*, Journal Computer Vision and Image Understanding, Vol.114, No.2, pp.274-285, 2010.
- [8] L.N. de Castro and J. Timmis, "Artificial immune systems: A new computational Intelligence Approach." Springer-Verlag, 1996.
- [9] D. Dasgupta, "Artificial immune systems and their applications," Springer-Verlag, 1999.
- [10] S.Gao, H.Dai, G.Yang, and Z.Tang, *A novel clonal selection algorithm and its application to travelling salesman problem*, IEICE Trans. Fundamentals, Vol.E90-A, pp.2318-2325, 2007.
- [11] T.Ichimura and S.Kamada, *Clustering and Retrieval Method of Immunological Memory Cell in Clonal Selection Algorithm*, Proc. of The 6th International conference on Soft Computing and Intelligent Systems and The 13th International Symposium on Advanced Intelligent Systems(SCIS-ISIS 2012), pp.1351-1356, 2012.
- [12] A.S. Perelson, *Immune network theory*, Immunological Review, Vol.110, pp.5-36, 1993.
- [13] J. Ng and S. Gong, *Learning Intrinsic Video Content Using Levenshtein Distance in Graph Partitioning*, Proc. of ECCV02, volume 4, pages 670684, Springer-Verlag, 2002.
- [14] T.Kohonen, *T. Self-Organizing Maps*, Springer Series in Information Sciences, Vol.30, Springer, 1995.
- [15] A.Rauber, D.Merkl, M.Dittenbach, *The growing hierarchical self-organizing map: exploratory analysis of high-dimensional data*, IEEE Transactions on Neural Networks, vol.13, pp.1331-1341, 2002.
- [16] R.A.Fisher, *Iris Dataset*, UCI Machine Learning Repository, <http://archive.ics.uci.edu/ml/datasets/Iris>, [online] 1936.
- [17] H.C.Wu, R.W.P.Luk, K.F.Wong, K.L.Kwok, *Interpreting TF-IDF term weights as making relevance decisions*, ACM Transactions on Information Systems, Vol.26, No.3, pp.137, 2008.



(b) With Levenshtein distance based similarity  
Fig. 9. Classification Result around Miyajima

Fig.9 shows the overview of classification result of subjective data by GHSOM. Fig.9(a) shows clusters by using only GPS, comments, and evaluation. The data without photos was classified by GHSOM, because the image data has a large amount of information and it is difficult for GHSOM to classify them as it is. On the contrary, Fig.9(b) is the classification result with the image which was processed by using the Levenshtein distance based similarity of images. In the paper, we prepare the images for 6 representative spots. Especially, the famous symbol 'Torii', a gateway at the entrance to a shrine, was drawn. The clusters in Fig.9(b) are divided into 3 groups; 'high similarity image with high TF-IDF and high evaluation', 'low similarity image with high TF-IDF and high evaluation', and 'low similarity image with low TF-IDF and low evaluation'. The data with 'high similarity image' represents the famous sightseeing spot. The data with 'high TF-IDF and high evaluation' and without 'high similarity image' means that the spot is including the cryptic tourist information. Therefore, the subjective data has scarcity value and will be a sightseeing spot potentially. Moreover, the classification of image data could distinguish the unknown spot in the field of famous landmark. In other words, the cryptic spots will be discovered in near landmark.