

# The Algorithm of Habitat Discovery in Bird Migration

Zhengzheng Wei<sup>1,\*</sup>, Zengzhen Shao<sup>1</sup>, Dong Chen<sup>1</sup>, and Yancong Li<sup>1</sup>

<sup>1</sup>School of Information Science and Engineering, Shandong Normal University, Jinan, 250300, P.R.China

**Abstract.** Bird migration has attracted an increasing attention. The study of habitats has played a vital role in the birds migratory. Previous researches, however, have encountered many problems, such as great limitations on research methods, low data utilization rate, statistics-focused and ineffective data processing and analysis methods. In this paper, the algorithm of habitat discovery is put forward by using computer's data-mining technology based on the spatio-temporal characteristics of bird-watching data. First the algorithm detects and eliminates duplicate data to guarantee data standardization. Then density-based clustering algorithms are used to identify habitats where birds gathered. Finally the habitats of birds migratory are discovered.

## 1 Introduction

All groups of birds are special in that a large number of species migrate annually between their breeding and non-breeding areas [1]. The migration of birds has a great impact on the environment and production life of human beings. Studying the migration of birds can help people prevent the spread of epidemics and maintain species diversity. Therefore the study of migration habitat is crucial for people to protect the birds and natural environment and maintain species diversity.

Various methods are adopted and developed to carry out research on migratory birds from different aspects by domestic and foreign researchers in order to understand the migration patterns of birds. Among these methods, the fixed-point investigation is the earliest bird migration research method. The most common and most popular method to study the bird migration is bird-banding[2]. It can be easily implemented and widely applied. But its monitoring cycle is long and the data recycling is quite complicated [3]. Next method is the satellite positioning method. Its accurately collected data can achieve continuous tracking to the individual bird. However it is not suitable for small birds with high cost and difficult popularization as well as limited amount of data. There are also some other methods including radar monitoring, sensitive geographical location and others. But they have low precision, difficulty in popularizing, limited data and other issues. In addition, the usage and analysis of the collected migration data of birds has also attracted the attention of researchers at home and abroad. It is analysed by the early bird data that only through the track point marked in GIS by biologists or the distribution points gotten by artificial statistics can migratory lands and migration routes be acquired [4]. In 2004, the Japanese scientist Shimazaki proposed to deal with bird flight data through using the method of ISODATA clustering [5]. In

this method the migratory state of birds is determined in accordance with their flight speed and further the migratory location is obtained. Nonetheless the migratory routes of birds still need to be manually marked and further the spatial location information is unable to be processed with this method. In 2010, Zhou Yuanchun and others found the birds gathering land by using the density-based hierarchical clustering algorithm to cluster bird GPS, the association of the aggregation rules and further migratory routes of birds by virtue of the Apriori algorithm or GSP algorithm [6]. However some defects that the changes of the habitats and migratory routes of birds cannot be found still exist due to small amount of data, less number of birds and short time span of the data. In 2012 Li Xueyan and others established a bird-watching database using China Birding Report and displayed the changes of bird distribution of recent years by GIS [7]. But it only relies on artificially statistical methods to mark bird discovery sites in the GIS, neither using the "quantitative observation" to analyse the data deeply, nor dealing with the problems of repeated sampling and uneven distribution for the bird-watching data.

From the above there are still a lot of problems that need to be solved in the above researches which are embodied in: 1) the collected data shows the following problems such as incomprehensive, imprecise, and limited amount; 2) there is no much work on data standardization of source data; 3) the amount of data for analysis and study is relatively small; 4) the hidden knowledge in the data fails to be dug out.

In this paper, the problem of traditional biology is abstracted as a computational problem. A feasible, efficient and general method is sought to solve the above problems, to achieve the effective treatment and utilization of bird-watching data, to find the habitats of migratory bird. This method can make up the shortcomings of the research of bird migration in China.

## 2 Data source

As the important supplementary information of traditional bird distributions, Chinese bird-watching data is comprehensive and reflects Chinese bird watching achievements accurately. These data comes from three aspects. The first is from the network such as Bird Report ([www.birdreport.cn](http://www.birdreport.cn)) [8] and China Bird Watching Network ([www.chinabirdnet.org](http://www.chinabirdnet.org)) [9] etc. The second comes from ornithological books and literature such as China Bird Report 2003-2007 [10], China Coastal Waterbird Census Report 2005-2011 [11] and *A Checklist and Distribution of the birds in Shandong* [12] etc. The third is provided by many ornithologists led by Prof. Sai Daojian. Total 189350 bird watching records have been verified by ornithologists that insure the authority of these data.

### 2.1 Data characteristics

The spatial-temporal information of Chinese bird-watching data records including species, date, location, number and observer in detail as shown in Fig. 1. The following “Number” records the number of birds observed in each bird watching.

BirdsID	BirdsEngName	BirdsLatinName	Address	Number	ActionTime	Observer
0257	Indian Jungle Nightjar	Caprimulgus indicus	Chongming Dongtan NR, Shanghai	6	2003/9/20	NR
0547	Black-crowned Night Heron	Nycticorax nycticorax	Fulan NR, Shenzhen, Guangdong	300	2003/1/26	HK
0343	Common Redshank	Tringa totanus	Huahu, Zaoze, Sichuan	30	2006/5/5	SY
0345	Common Greenshank	Tringa nebularia	Yellow River Delta NR, Dongying, Shandong	58	2006/4/15	SK
0096.1	Common Teal	Anas crecca	East Dongting Hu, Yueyang, Hunan	980	2006/2/12	LJY
0614	Brown Shrike	Lanius cristatus	Tianjin Normal University, Xiqing District, Tianjin	12	2015/5/20	NK
0518	Little Grebe	Tachybaptus ruficollis	Malanshan Forest Park, Hongshan District, Wuhan, Hubei	10	2015/6/14	HWJ
0542	Great Egret	Casmerodius albus	East Dongting Hu NR, Yueyang, Hunan	2	2015/12/12	XS
0349	Wood Sandpiper	Tringa glareola	Miyun Reservoir, Miyun, Beijing	12	2012/4/22	TT
0057	Whooper Swan	Cygnus cygnus	Yellow River Delta NR, Shandong	1000	2005/3/6	CHY
0058	Tundra Swan	Cygnus columbianus	Poyang Hu NR, Jiangxi	3497	2005/2/16	WWFC
0059	Swan Goose	Anser cygnoides	Shengqihu NR, Anhui	24211	2005/2/19	WWFC
0090	Spot-billed Duck	Anas platyrhynchos	Jiaocheng, Ningde, Fujian	150	2003/09/2	XY
0092	Northern Shoveler	Anas platyrhynchos	Fulan NR, Shenzhen, Guangdong	1000	2003/1/5	YY
0357	Great Knot	Calidris temminckii	Yalu Jiang NR, Dandong, Liaoning	1700	2006/5/13	BQ
0085	Gadwall	Anas strepera	Yellow River Wetland NR, Mengjin, Henan	108	2005/2/10	DQ

Fig. 1. Chinese Bird-Watching Data

### 2.2 Data Statistics

The 239350 Chinese bird-watching records involve 24 Orders, 100 Families, and 1230 Species accounting for 85.7% of China’s existing bird species [13]. The range of these records covers 34 provinces, municipalities and autonomous regions, including Hong Kong, Macao and Taiwan approximately 46 years ranging from 1970 to the present. Seeing from Chinn’s bird-watching records, the number of bird-watching records is more in the east than that in the west, and more in the south than that in the north. The proportions of bird-watching records also vary with years. The number of records from 2001 to 2016 accounts for more than half of the total. In addition, among the whole bird-watching records, there are 142078 records of migratory bird which accounts for 59.36% of the total. The remaining records are for non-migratory bird.

Although China’s bird-watching records are distributed in the species, time and space, the advantages of bird-watching data are obvious such as low cost, bulk information, easy access, various species coverage,

substantial data, long time span, high accuracy and convenient spatial-temporal data analysis [5].

## 3 Algorithm

### 3.1 Overview of Algorithm

Bird migration is a relatively long and complex process which reflects in spatial and temporal changes. The study of habitats has played a vital role in the birds migratory. In this paper Chinese bird-watching data is used as the data source. First the important attribute --"Number" (Fig 1) is introduced to solve the problems of repeated sampling and uneven sampling distribution based on the data characteristics of bird-watching records. The quality of the data will be improved. Then according to the temporal and spatial structure of the migration trajectory, the time and space attributions of the bird-watching records are dealt with separately. The potential information of the spatial -temporal data of bird is discovered. As a consequence, the quality of data mining will also be improved to some extent. This study will lead to more efficient utilization, processing and analysis of bird-watching data and provide new perspectives and new ideas for the study of bird migration in China.

Moreover migratory birds are employed as the study objects. This Algorithm is used to achieve two goals of the research: 1) to solve the problems of duplicate sampling and uneven distribution of sampling data; 2) to identify the habitats of migratory birds; These steps of the Algorithm are described in the following:

1) the bird-watching data is digitized and stored in the database as the form of GPS track point; 2) the temporal and spatial distance is calculated between the points to discover implicit duplicate data and the special points instead of the duplicated data are used to produce a standardized data set; 3) the preprocessed new trajectory point collection is clustered by the density-based clustering algorithm to obtain the high density area of the migratory activity that is used as a habitat during the migration.

The relevant definitions and steps of the Algorithm are described as follows in more detail.

### 3.2 Digitizing the Text into GPS Trajectory Points

Each bird-watching record contains its unique spatial-temporal data information. In order to better show the distribution and migration of birds each bird-watching record is abstracted into coordinate points with times before the mining and analysis steps. Each record corresponds to a coordinate point. This tracing point that can stand for an individual or a group is used to display the bird distribution and migration. The textual information of “bird-watching site” in each bird-watching record needs to be converted into the latitude and longitude coordinates of GPS in order to facilitate compare calculation and presentation. API interface of Baidu map with high accuracy is applied to the textual information of “bird-watching site” in each bird-watching

record. The coordinates of Baidu Map can be accurate after the 10 decimal point. By this Algorithm, the Migratory birds' corresponding coordinates are retrieved.

Definition1. Trajectory point of moving object is to describe the moving object sampling including three parts of latitude, longitude and timestamp and moving object identification which is expressed as  $p_i(e_j) = \langle x_i, y_i, t_i, e_j \rangle$ ,  $p_i(e_j) \in P(e_j)$ ,  $e_j \in E$ ,  $j \in [1, J]$ ,  $i \in [1, n]$  where  $(x_i, y_i)$  is the track point position component. For example,  $e_j = \text{Hirundo rustica}$ . A trajectory point is expressed as  $p_i(e_j) = \langle x_i, y_i, t_i, e_j \rangle$  where  $(x_i, y_i)$  are coordinates and  $t_i$  is a time stamp.  $p_i(e_j) \in P(e_j)$ ,  $j \in [1, J]$ ,  $i \in [1, n]$ .

Definition2. Migratory trajectory of migratory bird: a sequence of spatial locations with time stamps is called a migratory trajectory of birds. A migratory trajectory can be expressed as:  $P(e_j) = \{p_1(e_j), p_2(e_j), \dots, p_i(e_j), \dots, p_n(e_j)\}$  where  $p_i(e_j)$  is a sampling point of the trajectory and  $n$  is the number of sampling points.  $E_j$  is the moving object (event),  $e_j \in E$ ,  $E = \{e_1, e_2, \dots, e_j, \dots, e_J\}$  is the set of moving objects (events),  $j \in [1, J]$  and  $J$  is the moving object Event).

### 3.3 Selection of Feature Point

There is a considerable duplication in the original bird-watching data. If these duplicate data could not be removed, it would affect the quality of data mining. Therefore the original trajectory points of migratory bird should be pre-processed and removed duplicates before analysis of the data. The problem of uneven distribution of bird-watching records is solved initially, which lays the solid foundation for further analyzing of the data.

There are two types of duplicate data in bird-watching data. The first type is explicit duplicate data manifesting the same time and place, namely  $p_i(e_j) = p_{i+1}(e_j)$ . For such duplicates simple merge processing is needed. The second type is implicit duplicate data which is not easy to find. When a kind of bird repeatedly sampled over a continuous interval and a small regional range, its sample can be considered repeated. They manifest multiple sampling points of a bird species with similar temporal and spatial characteristics.

Definition3. Duplicate Points Set (DPS): the set of the tracking point of a certain type of moving object  $P'(e_j) = \{p_u(e_j), p_{u+1}(e_j), \dots, p_i(e_j), \dots, p_v(e_j)\}$ ,  $j \in [1, J]$ ,  $i \in [u, v]$ . If  $\text{Dist}(p_u, p_i) \leq 2 \theta r$ ,  $\text{Dist}(p_u, p_{v+1}) > 2 \theta r$ , and  $t_v - t_u < \theta t$ , the set  $P'(e_j)$  is called a Duplicate Points Set (DPS) marking as  $P'(e_j) = \{p_u'(e_j), p_{u+1}'(e_j), \dots, p_v'(e_j)\}$ . The sequence  $\{P_1'(e_j), P_2'(e_j), \dots, P_k'(e_j), \dots, P_K'(e_j)\}$  is called the set of DPSs where  $K$  is the number of DPS and  $k \in [1, K]$ . The tracking points of the DPS are implicit duplicate data. Fig 2 shows the DPSs.

Definition4. Feature Point (FP): In a DPS, a single point  $P'(e_j)$  that can be used to replace  $P'(e_j)$  is called a Feature Point (FP). The FP can be expressed as:  $sk(e_j) = \langle x_k, y_k, t_u, t_v, e_j \rangle$  where  $t_u$  is the starting time of the DPS,  $t_v$  is the ending time and  $(x_k, y_k)$  is the coordinate

of the FP,  $k \in [1, K]$ .

The steps of the selection of Feature Point are as follows:

- Step1: Assign the values of  $\theta r$  and  $\theta t$ ;
- Step2: A point  $p_i(e_j)$  is chose as a center arbitrarily in the set  $P(e_j) = \{p_1(e_j), p_2(e_j), \dots, p_i(e_j), \dots, p_n(e_j)\}$ ;
- Step3: Calculate the distances of the remaining points except the central point of the set  $P(e_j)$  to  $p_i(e_j)$ ;
- Step4: If  $\text{Distance}(p_{other}, p_i) \leq \theta r$  &&  $|t_{other} - t_i| < \theta t$ , then the point  $p_i(e_j)$  is added to the DPS  $P'(e_j)$ ;
- Step5: Output  $P'(e_j)$ ;
- Step6: The k-mediods algorithm (Han et al. 2012) is used in the DPSs  $P'(e_j)$ . The cluster of k-mediods algorithm is set to 1. The central point of cluster  $sk(e_j)$  is used as the feature point of  $p_i(e_j)$ ,  $sk(e_j) = \langle x_k, y_k, t_u, t_v, e_j \rangle$ ;
- Step7: Calculate the weighted average of points in  $P'(e_j)$ . The weight average of a point is the value of "Number" that is used as a new weight of  $sk(e_j)$ ;
- Step8: Repeat Step2-Step7 for each point in  $P'(e_j)$  until all FPs are output;
- Step9: Rearrange the feature points and the tracking points that are not in the DPSs. Then a set of new trajectory points about  $e_j$  is got,  $S(e_j) = \{s_1(e_j), s_2(e_j), \dots, s_k(e_j), \dots, s_m(e_j)\}$ ,  $m$  is the number of the new trajectory points;
- Step10: Output  $S(e_j)$ .

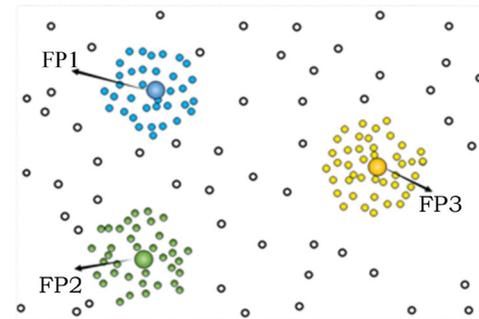


Fig. 2. Duplicate Points Sets (DPS)

### 3.4 The Discovery of Habitats

There are many stopover sites in the route of birds migrant which consist the habitants with wintering ground and breeding place [14][15]. Within the habitat, the number of bird populations is usually larger than that of the other regions [16][17]. That is a place that has larger amount and higher density of birds is much more important to birds themselves which is more likely to become the potential habitats for birds. Following this idea, the area where birds are densely distributed must be found after the pre-processing of the new trajectory points. A density-based clustering algorithm is employed in the new trajectory points  $S(e_j) = \{s_1(e_j), s_2(e_j), \dots, s_k(e_j), \dots, s_m(e_j)\}$  to find the high density areas for bird activity. And the value of "Number" is introduced to calculate the values of each cluster and outlier. Then high-density and high-volume areas are screened out as habitats for migratory bird. Fig 3 shows various types of clusters.

Definition5. Heat: The importance of a point or region is called the Heat Degree (HD). The “number of birds” is used as the weight of the point which is called the Heat Degree Point (HDP). If it has higher weight, the heat range will be larger. The greater the sum of HDPs is in a region, the greater the HD of the region will be. The steps of the discovery of habitats are as follows:

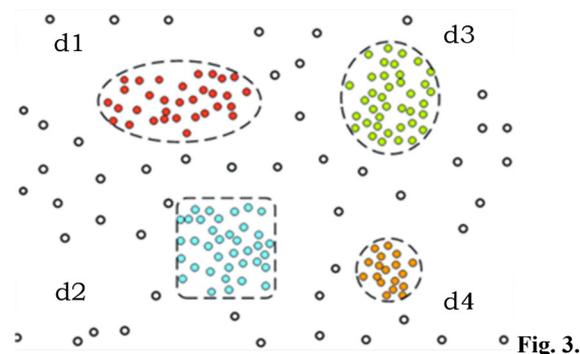
Step1: The DBSCAN algorithm (Han et al. 2012) is used in the new trajectory points sequence of  $e_j$ ,  $S(e_j) = \{s1(e_j), s2(e_j), \dots, sm(e_j)\}$  to get various shapes of clusters, denoted as  $C(e_j) = \{c1(e_j), c2(e_j), \dots, cl(e_j), \dots, cL(e_j)\}$ ,  $l \in [1, L]$ ,  $L$  is the number of the clusters;

Step2: The “number of birds” is used as the weight to calculate the HD of each cluster: the sum of multiplication of all points and the weights in a cluster  $c_l(e_j)$ . The HD of a cluster is called the HDC;

Step3: Calculate the HDP of each outlier point;

Step4: Arrange all HDPs and HDCs in ascending order. If  $HD > MinHeat$  (the value of  $MinHeat$  can be set by the user) the point or cluster will be output that is used as the habitat of migratory bird  $e_j$ :  $D(e_j) = \{d1(e_j), d2(e_j), \dots, dl(e_j), \dots, dL(e_j)\}$ ,  $l \in [1, L]$ ;

Step5: Repeat Step1 to Step4 for  $E = \{e1, e2, \dots, e_j, \dots, eJ\}$  until all habitats are output.



Clusters of Arbitrary Shape

### 3.5 Time complexity of Algorithm

There are two stages in this Algorithm. The time complexity of each stage will be separately analysed.

The k-medoids algorithm is used when selecting Feature Points (FP). The time complexity of the k-medoids algorithm is  $O(k(n-k)^2)$ . In this paper setting  $k=1$ ,  $n$  is the total number of trajectory points,  $S$  is the total number of points in the Duplicate Points Set (DPS) and  $S < n$ . The time complexity of algorithm is  $O((S/K)^2)$  when executing once and  $K$  is the number of Duplicate Points Sets (DPSs). This needs to be performed  $K$  times and  $K \ll n$ . Therefore the time complexity of this stage is  $O(S^2/K) < O(n^2)$ .

Then in the discovery of habitats the improved DBSCAN algorithm was executed once. Hence its time complexity is  $O(m^2)$  where  $m$  is the number of input trajectory points and  $m < n$  ( $n$  is the total number of trajectory points).

In conclusion the time complexity of the Algorithm is:  $O(S^2/K) + O(m^2) < O(n^3)$ . This shows that this

Algorithm can be completed in a relatively short time.

## 4 Experiment

### 4.1 Experimental environment

Bird-watching data is used for studying migratory bird in China thus the algorithm of habitat discovery is proposed. In the environment barn swallows (*Hirundo rustica*) is used as examples to explore their habitats. The feasibility and effectiveness of the algorithm is verified by comparing with the results of authoritative ornithological literatures *A Field Guide to the Birds of China*[18] and *A Checklist on the Classification and Distribution of the Birds of China*[19].

The operating environment of this experiment is windows7 operating system and C # language is used to write the algorithm. The software development environment is Microsoft Visual Studio 2010 and SQL server 2010.

### 4.2 Experiment One: Selection of Feature Point

The authors adopt the Feature Point (FP) to replace duplicate sampling points in the trajectory data of migratory bird.

Fig 4 shows partial data of a Duplicate Points Set (DPS) of *Hirundo rustica* from May 2, 2008 to May 9, 2008.

BirdsChinName	ActionTime	BirdsX	BirdsY
Hirundo rustica	2008/5/2 0:00:00	118.823466	31.97004
Hirundo rustica	2008/5/2 0:00:00	118.836165	31.072
Hirundo rustica	2008/5/2 0:00:00	118.781166	31.09214
Hirundo rustica	2008/5/3 0:00:00	118.836111	31.07212
Hirundo rustica	2008/5/3 0:00:00	118.856166	31.05534
Hirundo rustica	2008/5/3 0:00:00	118.800006	31.07215
Hirundo rustica	2008/5/4 0:00:00	118.8361	31.06214
Hirundo rustica	2008/5/4 0:00:00	118.743656	31.987115
Hirundo rustica	2008/5/4 0:00:00	118.801935	31.016251
Hirundo rustica	2008/5/4 0:00:00	118.685015	31.985155
Hirundo rustica	2008/5/5 0:00:00	118.713761	31.027857
Hirundo rustica	2008/5/5 0:00:00	118.809197	31.939081
Hirundo rustica	2008/5/5 0:00:00	118.78735	31.936139
Hirundo rustica	2008/5/6 0:00:00	118.779405	31.943003
Hirundo rustica	2008/5/6 0:00:00	118.781601	31.996915
Hirundo rustica	2008/5/6 0:00:00	118.757454	31.043939
Hirundo rustica	2008/5/6 0:00:00	118.760904	31.091918
Hirundo rustica	2008/5/8 0:00:00	118.698813	31.005734
Hirundo rustica	2008/5/8 0:00:00	118.752855	31.986915
Hirundo rustica	2008/5/9 0:00:00	118.763203	31.927214
Hirundo rustica	2008/5/9 0:00:00	118.843692	31.992015
Hirundo rustica	2008/5/9 0:00:00	118.794249	31.952808
Hirundo rustica	2008/5/9 0:00:00	118.765865	31.956729

Fig. 4. Partial data of a DPS (*Hirundo rustica*)

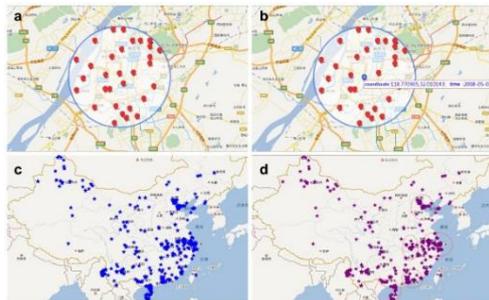
In the light of the bird-watching data and spherical distances based on latitude and longitude coordinates the following parameters are shown in Table1. After experimental tests the clustering results are best when  $\theta_r=6$  km and  $\theta_t=10$  days.

Table1. Parameters of DPS

NO.	Event E	$\theta_r(\text{km})$	$\theta_t(\text{day})$
1	$e_j$	6	10
2	$e_j$	6	15

3	ej	10	10
4	ej	10	15
5	ej	3	10
6	ej	3	15

When  $e_j = \text{Hirundo rustica}$  the trajectory points in Fig 4 (on May 2, 2008 to May 9, 2008, near Nanjing, Jiangsu Province) are marked on the map shown in Fig 5a. After clustering the center of the cluster – FP is gotten (Fig 5b) which is used to replace all DPSs by this method. After data pre-processing, the trajectory points become fewer (shown in Fig 5d) than that of the original (shown in Fig 5c).

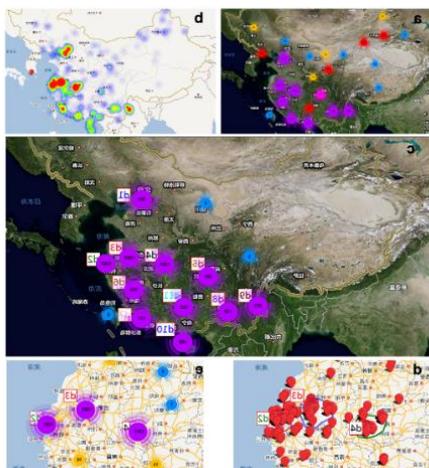


**Fig.5.** Data pre-processing of the Trajectory Points (*Hirundo rustica*).  
 a Duplicate Points Set. b Selecting the Feature Point. c The Original Trajectory Points d New Trajectory Points after Data Pre-processing

### 4.3 Experiment Two: The discovery of habitats

For the new trajectory points  $S(e_j)$  after preprocessing the density-based clustering algorithm and Heat Degree (HD) formula is applied to get habitats for avian migration.

When  $e_j = \text{Hirundo rustica}$  25 clusters  $C(\text{Hirundo rustica})$  can be gotten through DBSCAN clustering that is shown on the map Fig 6a. The HD of each point is then calculated and shown on the heat map Fig 6b. Finally 11 habitats are output.  $D(\text{Hirundo rustica}) = \{d_1, d_2, d_3, d_4, d_5, d_6, d_7, d_8, d_9, d_{10}, d_{11}\}$ , their  $HD = \{369, 1287, 1885, 1221, 287, 155, 177, 159, 231, 189, 122\}$  (Fig 6c).



**Fig .6.** The Discovery of Habitats (*Hirundo rustica*)  
 a The Clusters of Trajectory Points after DBSCAN. b The Distribution Diagram of the Heat Degree c Habitats d Trajectory Points within the Habitat. e Part of the Habitat

## 5 Discussion

### 5.1 Social and scientific value

This Algorithm is general, practical and convenient. Ideally it can be applied to all bird-watching data of migratory bird in China. However the algorithm depends on a certain scale of bird-watching data. If the quantity of data is greater, the results of data mining may be more precise. If the data sample is too small, the accuracy of the algorithm will decline and the distribution and migration of migratory bird cannot be reflected truly. These results are confirmed by the experiments. Therefore it is important to collect bird-watching data constantly and adequately data pre-processing is necessary before data analysis and mining. On the one hand, the algorithm solves the repeated sampling of data sets and ensures the accuracy of data mining. On the other hand, massive and redundant data is compressed by this process. Thus the efficiency of data analysis is also improved.

### 5.2 Limitations and shortcomings

1) Bird-watching data cannot track individual bird. Therefore it is not enough to just rely on bird-watching data itself to verify the accuracy of the results. 2) There are spatial and temporal discontinuities in bird-watching data. Hence it is difficult to solely rely on such data to analysis changes of bird migration's habitats over the years. Moreover this may insert a negative impact on the predicted results of bird migration. 3) The uneven distribution of sampling still exists which prevents deeper data analysis and mining activities.

### 5.3 Future directions

First bird-watching data will be collected continuously. Additional migratory bird data such as satellite-tracking data and bird-banding data will be gradually introduced to supplement further mining and analyses. Second meteorological data will be added to our study. The authors will further study the effects of climate and environment on bird migratory habitats.

## Conclusion

Bird watching activity has been developing rapidly in China in recent years. This activity will help people to understand the distributions of birds and population dynamics etc. This paper analyses the problems of bird migration in-depth from the perspective of mining data. Based on Chinese bird-watching records, the algorithm of habitat discovery is proposed and used for the selection

and discovery of habitats during the migratory process of birds. Taking *Hirundo rustica* as example, maps and GIS demonstrate the feasibility of the algorithm. The time complexity of the algorithm is small, resulting in its high efficiency. The migration routes and habitats of the birds derived by this work are compared with that of the authoritative ornithological literature which shows more accurate and real results and verifies the accuracy of the algorithm.

## References

1. Bock WJ. Presidential Address: Three Centuries of International Ornithology. *Acta Zoologica Sinica*. 2004, 50(6): 880-912
2. Wang Y, Zhang ZW, Zheng GM, Li JQ, Xu JL, Ma ZJ. Ornithological Research: Review of the Past Twenty Years and Future Perspectives in China. *Biodiversity Science*. 2012, 20: 119-137
3. Ma ZJ. The Research Methods and Study Advances of Bird Migration. *Bulletin of Biology*. 2009, 44:5-9
4. Bousquet O, Luxburg UV, Rätsch G. *Advanced Lectures on Machine Learning*. Springer-Verlag, 2003.
5. Shimazaki H, Tamura M, Higuchi H. Migration Routes and Important Stopover Sites of Endangered Oriental White Storks (*Ciconia boyciana*) as Revealed by Satellite Tracking. *Revista Enfermería Del Trabajo*. 2004, 3:39-43
6. Zhou YC, Tang MJ, Cui P. Research and Implementation of Data Mining Algorithm for Bird Migratory Behavior in Qinghai Lake. *E-Science Technology and Application*. 2010, 1:38-50
7. Li XY, Liang LY, Gong P, Liu Y, Liang FF. Revealing Bird Distribution Changes of Bird Watching in China. *Chinese Science Bulletin*. 2012, 57: 2956-2963.
8. Bird Report. <http://www.birdreport.cn/>. 2016. Accessed 23 Apr 2016.
9. Bird Watching Network of China. <http://www.chinabirdnet.org/indexc.html>. 2016. Accessed 1 Apr 2016.
10. China Ornithological Society. *China Bird Report 2007 (in Chinese)*. Beijing: China Ornithological Society, 2008
11. China Coastal Waterbird Census Team. *China Coastal Waterbird Census Report (in Chinese)*. Hong Kong: Bird Watching Society, 2011
12. Sai DJ, Sun YG. *A Checklist and Distribution of the birds in Shandong*. Beijing: Science Press, 2013
13. Liu Y, Wei Q, Dong L. Updated New Bird Records in China Recently. *Chinese Journal of Zoology*. 2013, 48:750-758
14. Ma ZJ, Li B, Chen JK. Study on the Utilization of Stopover Sites and Migration Strategies of Migratory bird. *Acta Ecologica Sinica*. 2005, 25:1404-1412.
15. Ma ZJ, Li B, Chen JK. Physiological Ecology of Migratory bird during the Stopover Periods. *Acta Ecologica Sinica*. 2005, 25:3067-3075.
16. Yang Y, Wen JB, Hu DF. A Review on Avian Habitat Research. *Scientia Silvae Sinicae*. 2011,47:172-180.
17. Munster VJ, Baas C, Lexmond P, Wallensten A, Fransson T, Rimmelzwaan GF, et al. Temporal, and Species Variation in Prevalence of Influenza A Viruses in Wild Migratory bird. *Plos Pathogens*. 2007, 3:630-638.
18. MacKinnon J R, Phillipps K and He F. *The Field Guide to the Birds of China*. Changsha: Hunan Education Press, 2010.
19. Zheng GM. *The Checklist of the Classification and Distribution of the Birds of China (Second Edition)*. Beijing: Science Press, 2011.