

Exploration of Interesting Dense Regions on Spatial Data

Double-blind submission

ABSTRACT

The large size of spatial data hinders its effective analysis for insight discovery. Analysts require to obtain few high quality options (i.e., “highlights”) to have a more focused investigation. In this paper, we define, formalize and analyze Interesting Dense Regions (IDRs) which capture implicit preferences of analysts in order to automatically find highlights. Our approach involves a polygon-based abstraction layer for captured preferences. Using these interesting dense regions, we highlight points to guide analysts in the analysis process. We discuss the efficiency and effectiveness of our approach through a realistic example and a few experiments on Airbnb and Yelp datasets.

KEYWORDS

Spatial data, implicit feedback, dense regions, spatial highlights.

ACM Reference Format:

Double-blind submission. 2018. Exploration of Interesting Dense Regions on Spatial Data. In *Proceedings of ACM International Conference on Information and Knowledge Management (CIKM’18)*. ACM, New York, NY, USA, Article 4, 4 pages. https://doi.org/10.475/XXX_X

1 INTRODUCTION

Spatial data are often voluminous. Hence the focus in the literature of spatial data analysis is often on “efficiency”, i.e., enabling fluid means of navigation in spatial data in order to facilitate the exploratory analysis. The common approach is to design pre-computed indexes which enable efficient retrieval of spatial data (e.g., [13, 16]). However, there has been fewer attention to the “value” derived from spatial data. Despite the huge progress on the efficiency front, an analyst may easily get lost in the plethora of geographical points due to two following reasons.

- In an exploratory context, the analyst doesn’t know apriori what to investigate next.
- Moreover, she may easily get distracted and miss interesting points by visual clutter caused by huge point overlaps.

The main drawback of the traditional analysis models (e.g., [12, 14, 15], to name a few) is that the analyst has a *passive role* in the process. In other words, the analyst’s feedback (i.e., her likes and dislikes) is ignored and only the input query (i.e., her explicit request) is served. In case feedback is incorporated, the process can be more directed towards analyst’s interests where her partial needs can be served earlier in the process. In this paper, we advocate for a “guidance layer” on top of the raw visualization of spatial data to enable analysts know “*what to see next*”. This guidance should

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

CIKM’18, October 2018, Lingotto, Turin Italy

© 2018 Copyright held by the owner/author(s).

ACM ISBN XXX-XXX-XX-XXX/18/10/18.

https://doi.org/10.475/XXX_X

be a function of analyst feedback: the system should recommend options similar to what the analyst has already appreciated.

Considering this context, we propose in this paper an approach to explore Interesting Dense Regions (IDRs) whose aim is to capture and analyze implicit feedback of analysts in spatial data analysis. Then a set of highlights is obtained from IDRs to guide analysts. Without loss of generality, we focus on “mouse moves” as the implicit feedback received from the analyst. Mouse moves are the most common way that analysts interact with geographical maps [7]. However, our approach can be easily extended to other types of inputs such gaze tracking [2] and leap motions.

2 INTERESTING DENSE REGIONS

An Interesting Dense Region (IDR) is a spatial region with high probability of having spatial points with potential interest to the analyst. An IDR is captured and defined by processing user feedbacks. Unlike the literature which mainly focuses on explicit reflections, we investigate on implicit feedback. There are different ways to capture implicit feedbacks.

- During spatial data analysis, it is often the case that analysts look at some regions of interest but forget to provide an explicit feedback. We call this latent signal, gaze. It is shown in [2] that gaze has a strong correlation with “user attention”. The signal can be captured by tracking eye movements.
- To address privacy issues of gaze exploitation, we consider an alternative option, i.e., tracking the mouse cursor. It is shown in [3] that mouse gestures have a strong correlation with “user engagement”. Intuitively, a point receives a positive feedback if the cursor moves around it frequently.

The objective of discovering IDRs is to obtain the preferences of analysts which have never been expressed explicitly. For such discovery, we are inspired by two observations.

■ **Observation 1.** We believe that a region appeals more interesting to the analyst if it is denser, i.e., the analyst moves her mouse in that region several times.

■ **Observation 2.** It is possible that the analyst moves her mouse everywhere in the map. This should not signify that everywhere in the map has the same significance.

We consider two different layers on a geographical map: “spatial layer” and “interaction layer”. The spatial layer contains points from a spatial database \mathcal{P} . The interaction layer contains mouse move points \mathcal{M} . Our proposed approach mines a set of mouse move points \mathcal{M} in the interaction layer to discover one or several IDRs, in which most analyst’s interactions occur. Then it matches the spatial points \mathcal{P} with IDRs in order to find points inside each region. We employ a Quadtree index as a mediator between the spatial and interaction layers in order to reduce the number of matchings. The attributes of resulting points will be exploited to update analyst’s feedback vector F . The updated vector F will then be used to find k highlights (where k is specified by the analyst). These steps ensure that the final highlights reflect analyst’s implicit interests.

Algorithm 1: Information Highlighting

```

Input: Current time  $t_c$ , mouse move points  $\mathcal{M}$ ,  $k$ 
Output:  $k$  highlights  $\mathcal{H}$ 
1  $S \leftarrow \emptyset$ 
2  $g \leftarrow \text{number of time segments}$ 
3 for  $i \in [0, g]$  do
4    $\mathcal{M}_i \leftarrow \{m = \langle x, y, t \rangle | (\frac{t_c}{g} \times i) \leq t \leq (\frac{t_c}{g} \times (i + 1))\}$ 
5    $C_i \leftarrow \text{mine\_clusters}(\mathcal{M}_i)$ 
6    $O_i \leftarrow \text{find\_polygons}(C_i)$ 
7 end
8 for  $O_i, O_j$  where  $i, j \in [0, g]$  and  $i \neq j$  do
9    $S.append(\text{intersect}(O_i, O_j))$ 
9  $\mathcal{H} \leftarrow \text{discover\_highlights}(S, k)$ 
10 return  $\mathcal{H}$ 

```

Algorithm 1 summarizes our approach for mining IDRs. We add points to \mathcal{M} only every 200ms to prevent adding redundant points. Following Observation 1 and in order to mine the recurring behavior of the analyst, the algorithm begins by partitioning the set \mathcal{M} into g fixed-length consecutive segments \mathcal{M}_0 to \mathcal{M}_g . The first segment starts at time zero (where the system started), and the last segment ends at t_c , i.e., the current time. Following Observation 2, we then find dense clusters in each segment of \mathcal{M} using a variant of DB-SCAN approach [9]. Finally, we return intersections among those clusters as IDRs.

For clustering points in each time segment (i.e., line 5 of Algorithm 1), we use ST-DBSCAN, a space-aware variant of DB-SCAN for clustering points based on density [6]. For each subset of mouse move points \mathcal{M}_i , $i \in [0, g]$, ST-DBSCAN begins with a random point $m_0 \in \mathcal{M}_i$ and collects all density-reachable points from m_0 using a distance metric. As mouse move points occur in the 2-dimensional pixel space (i.e., the display), we choose euclidean distance as our metric. If m_0 turns out to be a core object, a cluster will be generated. Otherwise, if m_0 is a border object, no point is density-reachable from m_0 and the algorithm picks another random point in \mathcal{M}_i . The process is repeated until all of the points have been processed.

Once clusters are obtained for all subsets of \mathcal{M} , we find their intersections to locate recurring regions (line 6). To obtain intersections, we need to clearly define the spatial boundaries of each cluster. Hence for each cluster, we discover its corresponding polygon that covers the points inside. For this aim, we employ Quickhull algorithm, a quicksort-style method which computes the convex hull for a given set of points in a 2D plane [4].

There exist several approaches to infer a spatial region for a given set of points [1, 5, 8, 10, 11]. The common approach is to cluster points in form of concave and convex polygons. In case a concave polygon is constructed, the “dents” of such a polygon may entail points which are not necessarily in \mathcal{M} . In IDR algorithm, however, we adapt Quickhull [4], due its simplicity, efficiency and it’s natural implementation of convex polygons.

Once the IDRs are discovered, we employ GEOGUIDE, an algorithm proposed in [15], to highlight k points among the set of points inside IDRs (line 9). We seek two properties in k highlights,

i.e., *similarity* and *diversity*. First, highlights should be in the same direction of the analyst’s implicit feedback, hence similar to the vector F . Second, highlighted points should also represent distinct directions so that the analyst can observe different aspects of data and decide based on the big picture. The similarity between points is captured using Cosine on their attributes. Diversity operates on geographical coordinates of points where farther points contribute more to diversity. The distance between two points is computed using Haversine measure [18].

3 CASE STUDY

We discuss a real-world example to show the functionality of our approach in practice.

Example. Lucas is planning to spend few days in Paris, France. His appreciation of French culture makes him interested in new experiences in the city. He decides to rent a home-stay from Airbnb website¹. He likes to discover the city, hence he is open to any type of lodging in any region with an interest to stay in the city center. The website returns 4000 different locations. As he has no other preferences, an exhaustive investigation needs scanning each location independently which is nearly infeasible. While he is scanning few first options, he shows interest in the region of “Champ de Mars” (near Eiffel Tower), but he forgets or doesn’t feel necessary to click a point there. By discovering IDRs on his mouse moves over the home-stays in Paris, our system can quickly detect his interest in the region and short-list a small subset of locations (i.e., highlights) accordingly to be recommended to Lucas.

We follow the above example to describe how IDRs are constructed in action. Figure 1 shows Lucas’ steps to explore home-stays in Paris. Figure 1.A shows his mouse movements in different time stages. In this example, we consider $g = 3$ and capture Lucas’ feedback in three different time segments (progressing from Figures 1.B to 1.D). It shows that Lucas started his search around Eiffel Tower and Arc de Triomphe (Figure 1.B) and gradually showed interest in south (Figure 1.C) and north (Figure 1.D) as well. All intersections between those clusters are discovered (hatching regions in Figure 1.E) which will constitute the set of IDRs (Figure 1.F), i.e., IDR1 to IDR4.

In the Airbnb dataset, points are home-stays. Figure 2 shows these points with their nightly price for IDRs. We observe that there exist many matching points with IDR3 and absolutely no matching point for IDR2. For IDR4, although there exist many home-stays below the region, we never check their containment, as they belong to a Quadtree cell which doesn’t intersect with the IDR.

4 EXPERIMENTS

We discuss two sets of experiments. Our first set is on the usability of our approach. Then we focus more on discovering IDRs and present few statistics and insights for them. In the interest of space, we only present a glimpse of our experiments here. More will be presented in an extended report which we will make available upon acceptance.

First off, we validate the “usability” of our approach. For this aim, we design a user study with some participants who are all

¹<http://www.airbnb.com>

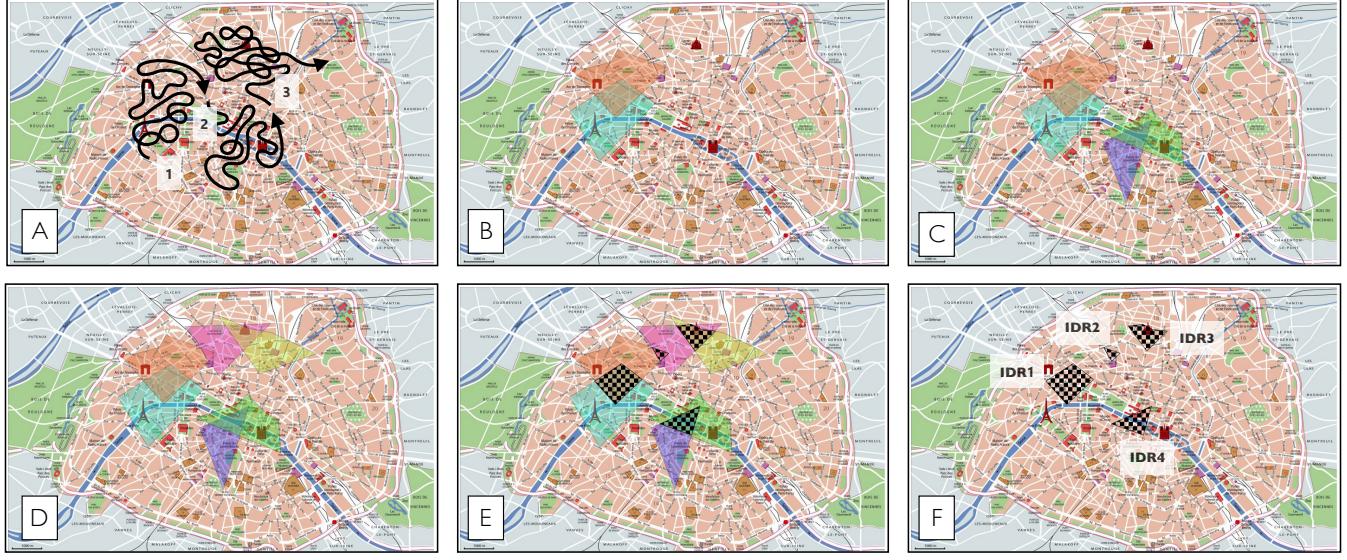


Figure 1: The process of finding IDR over Paris home-stays.

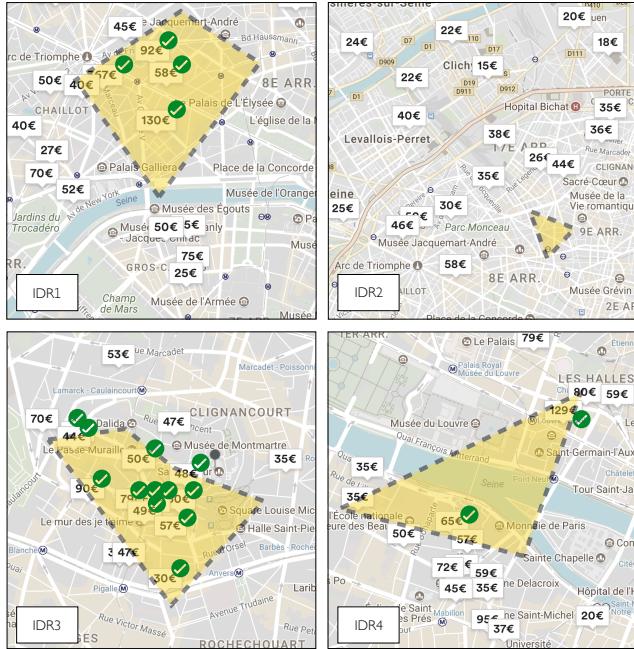


Figure 2: Matching points for IDR.

students of Computer Science. Some of them are “novice” users who don’t know the location under investigation, and some are “experts.” Participants should fulfill a task. For each participant, we report a variant of time-to-insight measure, i.e., how long the participants interact with the tool before fulfilling the task. Evidently, less number of interactions is preferred as it means that the participant can reach insights faster.

On the Airbnb dataset of Paris with 1,000 points, we define two different tasks: T1: “finding a point in a requested location” (e.g., find a home-stay in the “Champ de Mars” area), and T2: “finding a point with a requested profile” (e.g., find a cheap home-stay.) Participants may also begin their navigation either from I1: “close to the goal” or I2: “far from the goal”.

Table 1: Interactions of “novice” and “expert” participants (in seconds)

	T1/I1	T2/I1	T1/I2	T2/I2
Novices	1.99	2.38	2.00	2.48
Experts	1.72	2.09	1.70	2.14

Table 1 shows the results. We observe that on average it takes 2.067 seconds to achieve defined goal. This shows that implicit feedback capturing is an effective mechanism which helps analysts to reach their goals much faster. Expert participants need 0.35 seconds less time on average. Interestingly, starting points, i.e., I1 and I2, do not have a huge impact on number of steps. It is potentially due to the diversity component which provides distinct options and can quickly guide analyst towards their region of interest. We also observe that the task T2 is an easier task than T1. This is potentially due to where the analyst can request options similar to what she has already observed and greedily move to her preferred regions.

In the second part of our experiments, we employ two different datasets, i.e., Airbnb² and Yelp³. We pick a similar subset from both datasets, i.e., home-stays and restaurants in Paris city, respectively. We consider four different sizes of those datasets, i.e., 100, 1000, 2000 and 4000 points, respectively. For each size of the datasets, we

² <http://insideairbnb.com/get-the-data.html>

³ <https://www.yelp.com/dataset>

manually perform 20 sessions, and then we present the results as the average of sessions.

We limit each session to 2 minutes where we seek for interesting points in the datasets. We capture the following information in each session:

- The number of regions created from the mouse moves during the session;
- The number of generated IDR (intersection of regions);
- The number of points from the dataset presented in each IDR;
- The coverage of points (in the dataset) with IDR collectively.

Tables 2 and 3 show the result for Airbnb and Yelp, respectively. In Table 2, we observe that the number of regions decreases when the number of points increases. On average, 10 regions are constructed per session. The average number of points presented in IDR is 24.97. It shows that our approach is able to highlight at least 8.05% of points from the dataset, on average.

Table 2: IDR statistics on Airbnb dataset

# points	# regions	# IDR	# points in IDR	% points
100	11.35	10.05	29.40	29.40%
1000	10.75	6.75	11.70	1.17%
2000	7.37	3.63	5.63	0.003%
4000	10.30	10.15	53.15	1.33%
average	9.94	7.64	25.97	8.05%

More uniform results are observed in Table 3, i.e., for Yelp dataset vis-à-vis Airbnb. The average number of generated regions reaches 12.75 per session. Also, the number of regions decreases by increasing the number of points. The same happens for IDR, where we obtain an average of 8.9 IDR generated per session. The number of points presented in IDR is on average 108.65 and it represents on average 13.11% of points highlighted from the dataset.

Table 3: IDR statistics on Yelp dataset

# points	# regions	# IDR	# points in IDR	% points
100	14.90	7.55	28.30	28.30%
1000	13.90	10.00	149.55	14.96%
2000	11.05	9.80	111.05	5.55%
4000	10.45	8.55	145.7	3.64%
average	12.57	8.97	108.65	13.11%

5 CONCLUSION

In this paper, we present an approach to explore Interesting Dense Regions (IDRs) using implicit feedback in order to detect analyst latent preferences. The implicit feedbacks are captured from mouse moves of analysts over the geographical map while analyzing spatial data. We formalize a novel polygon-based mining algorithm which returns few highlights in line with analyst's implicit preferences. The highlights enable analysts to focus on what matters the most and prevent information overload.

We consider various directions of future work for this work. First, we are interested to incorporate an "explainability" component

which can describe causalities behind preferences. For instance, we are interested to find seasonal patterns to see why the preferences of analysts change from place to place during various seasons of the year. Another direction is to incorporate "Query by Visualization" approaches, where analysts can specify their intents alongside their implicit preferences, directly on the map [17].

REFERENCES

- [1] Avi Arampatzis, Marc van Kreveld, Iris Reinbacher, Christopher B. Jones, Subodh Vaid, Paul Clough, Hideo Joho, and Mark Sanderson. 2006. Web-based delineation of imprecise regions. *Computers, Environment and Urban Systems* 30, 4 (2006), 436 – 459. <https://doi.org/10.1016/j.compenvurbsys.2005.08.001> Geographic Information Retrieval (GIR).
- [2] Ioannis Arapakis, Mounia Lalmas, B Barla Cambazoglu, Mari-Carmen Marcos, and Joemon M Jose. 2014. User engagement in online news: Under the scope of sentiment, interest, affect, and gaze. *Journal of the Association for Information Science and Technology* 65, 10 (2014), 1988–2005.
- [3] Ioannis Arapakis, Mounia Lalmas, and George Valkanas. 2014. Understanding within-content engagement through pattern analysis of mouse gestures. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*. ACM, 1439–1448.
- [4] C. Bradford Barber, David P. Dobkin, and Hannu Huhdanpaa. 1996. The Quickhull Algorithm for Convex Hulls. *ACM Trans. Math. Softw.* 22, 4 (Dec. 1996), 469–483. <https://doi.org/10.1145/235815.235821>
- [5] Michael Bevis and Jean-Luc Chatelain. 1989. Locating a point on a spherical surface relative to a spherical polygon of arbitrary shape. *Mathematical Geology* 21, 8 (01 Oct 1989), 811–828. <https://doi.org/10.1007/BF00894449>
- [6] Derya Birant and Alp Kut. 2007. ST-DBSCAN: An Algorithm for Clustering Spatial-Temporal Data. *Data Knowl. Eng.* 60, 1 (Jan. 2007), 208–221. <https://doi.org/10.1016/j.datnak.2006.01.013>
- [7] Mon Chu Chen, John R. Anderson, and Myeong Ho Sohn. 2001. What Can a Mouse Cursor Tell Us More?: Correlation of Eye/Mouse Movements on Web Browsing. In *CHI '01 Extended Abstracts on Human Factors in Computing Systems (CHI EA '01)*. ACM, New York, NY, USA, 281–282. <https://doi.org/10.1145/634067.634234>
- [8] Matt Duckham, Lars Kulik, Mike Worboys, and Antony Galton. 2008. Efficient generation of simple polygons for characterizing the shape of a set of points in the plane. *Pattern Recognition* 41, 10 (2008), 3224 – 3236. <https://doi.org/10.1016/j.patcog.2008.03.023>
- [9] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. 1996. A Density-based Algorithm for Discovering Clusters a Density-based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD'96)*. AAAI Press, 226–231. <http://dl.acm.org/citation.cfm?id=3001460.3001507>
- [10] M.J. Fadili, M. Melkemi, and A. ElMoataz. 2004. Non-convex onion-peeling using a shape hull algorithm. *Pattern Recognition Letters* 25, 14 (2004), 1577 – 1585. <https://doi.org/10.1016/j.patrec.2004.05.015>
- [11] Antony Galton and Matt Duckham. 2006. What Is the Region Occupied by a Set of Points?. In *Geographic Information Science*, Martin Raubal, Harvey J. Miller, Andrew U. Frank, and Michael F. Goodchild (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 81–98.
- [12] Niranjan Kamat, Prasanth Jayachandran, Karthik Tunga, and Arnab Nandi. 2014. Distributed and interactive cube exploration. In *Data Engineering (ICDE), 2014 IEEE 30th International Conference on*. IEEE, 472–483.
- [13] Lauro Lins, James T Klosowski, and Carlos Scheidegger. 2013. Nanocubes for real-time exploration of spatiotemporal datasets. *IEEE Transactions on Visualization and Computer Graphics* 19, 12 (2013), 2456–2465.
- [14] Behrooz Omidvar-Tehrani, Sihem Amer-Yahia, and Alexandre Termier. 2015. Interactive User Group Analysis. In *CIKM*. ACM, 403–412. <https://doi.org/10.1145/2806416.2806519>
- [15] Behrooz Omidvar-Tehrani, Plácido A Souza Neto, Felipe M Freire Pontes, and Francisco Bento. 2017. GeoGuide: An Interactive Guidance Approach for Spatial Data. In *Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData), 2017 IEEE International Conference on*. IEEE, 1112–1117.
- [16] Cícero AL Pahins, Sean A Stephens, Carlos Scheidegger, and Joao LD Comba. 2017. Hashedcubes: Simple, low memory, real-time visual exploration of big data. *IEEE transactions on visualization and computer graphics* 23, 1 (2017), 671–680.
- [17] Tarique Siddiqui, Albert Kim, John Lee, Karrie Karahios, and Aditya Parameswaran. 2016. Effortless data exploration with zenvisage: an expressive and interactive visual analytics system. *Proceedings of the VLDB Endowment* 10, 4 (2016), 457–468.
- [18] Glen Van Brummelen. 2012. *Heavenly mathematics: The forgotten art of spherical trigonometry*. Princeton University Press.