



東南大學
SOUTHEAST UNIVERSITY

Efficient Similar Region Search with Deep Metric Learning

Yiding Liu, Kaiqi Zhao, Gao Cong Nanyang Technological University

KDD '18: The 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining

张斌杰
2021.04.15

目录



1

问题背景

2

区域相似性

3

搜索算法

4

实验

5

讨论

01

问题背景



问题背景

data in city [1] :

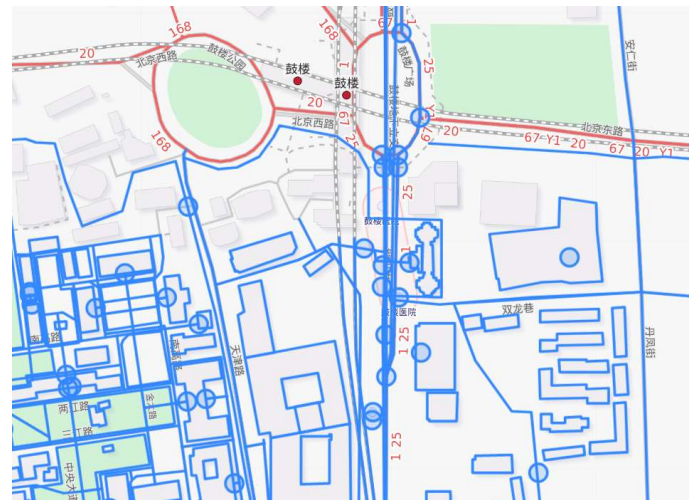
- POI
- Road
- Mobility Flow
-

applications of urban computing:

- POI recommendation
- Chain-store placement recommendation
-



POI data



Road data



Mobility Flow data[2]

- [1] Xu, Y., Shen, Y., Zhu, Y., & Yu, J. (2020). Ar2Net: An attentive neural approach for business location selection with satellite data and urban data. ACM Transactions on Knowledge Discovery from Data, 14(2). <https://doi.org/10.1145/3372406>
- [2] Jenkins, P., Farag, A., Wang, S., & Li, Z. (2019). Unsupervised representation learning of spatial data via multimodal embedding. International Conference on Information and Knowledge Management, Proceedings, 1993–2002. <https://doi.org/10.1145/3357384.3358001>

问题背景



Usage of Similar Area Search:

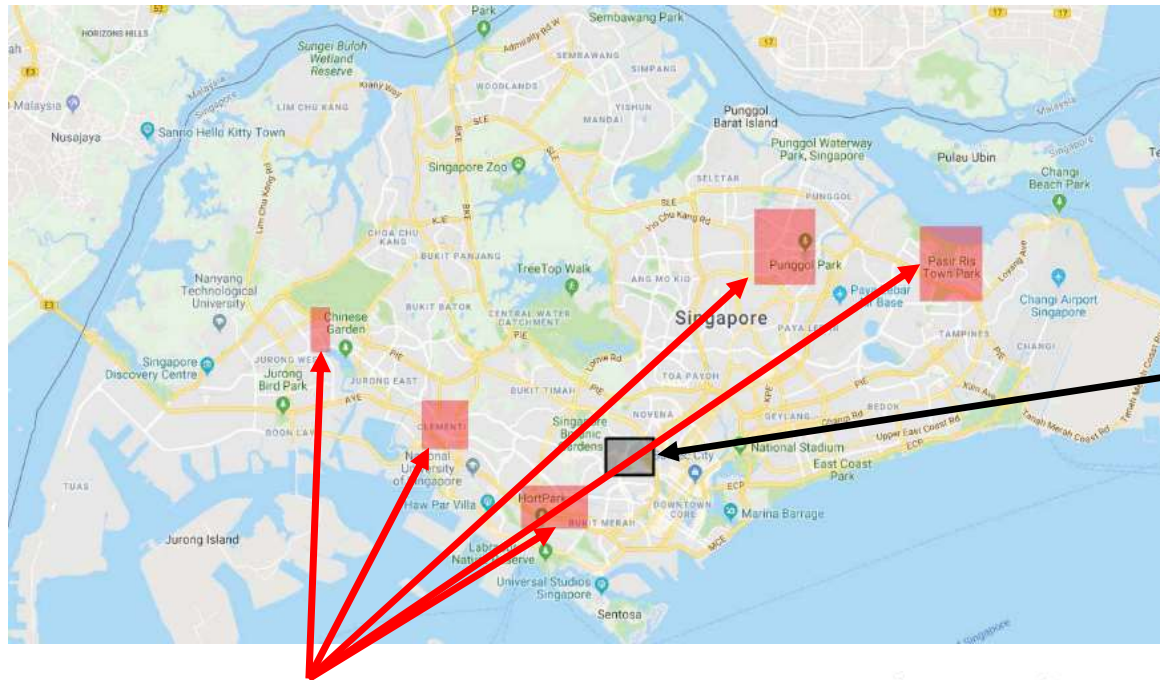
- Business site selection
- Improving location-based services. (POI recommendation)
- Explore unfamiliar cities or regions
-



问题背景



What is Similar Area Search:



geographical space P

query region R_q

a similarity metric $\text{sim}(\cdot)$

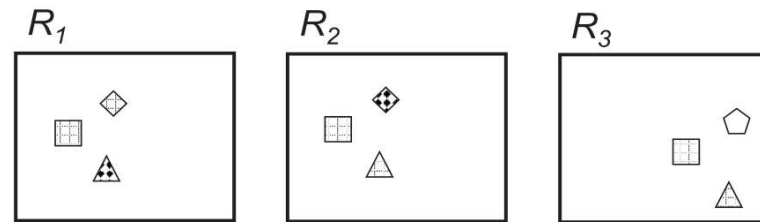
retrieve a set of N regions as \mathcal{R} .

$$\text{sim}(R_q, R_i) \geq \text{sim}(R_q, R_j), \forall R_i \in \mathcal{R}, \forall R_j \notin \mathcal{R}.$$

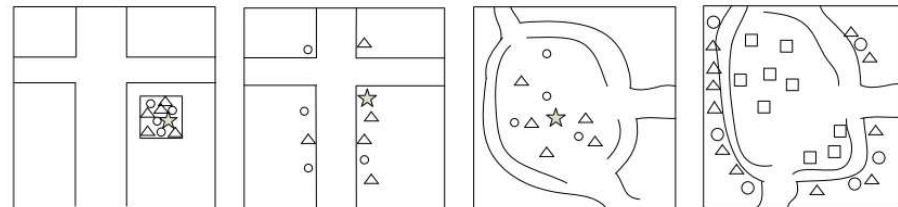
Challenges of Similar Area Search:

- How to model region similarity

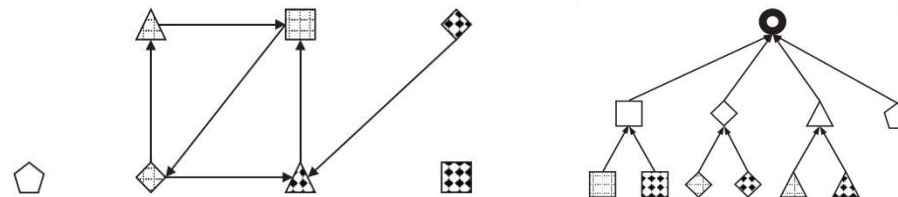
categories of POIs inside[1]



distribution of POIs inside [2]



relations of POIs/categories inside [1]



[1] Jin, X., Lee, D., Oh, B., Lee, K. H., Lee, S., & Chen, L. (2019). Learning region similarity over spatial knowledge graphs with hierarchical types and semantic relations. International Conference on Information and Knowledge Management, Proceedings, 669–678. <https://doi.org/10.1145/3357384.3358008>

[2] Sheng, C., Zheng, Y., Hsu, W., Lee, M. L., & Xie, X. (2010). Answering Top-k Similar Region Queries. In International Conference on Database Systems for Advanced Applications: Vol. 5981 LNCS (Issue PART 1, pp. 186–201). https://doi.org/10.1007/978-3-642-12026-8_16

问题背景



Challenges of Similar Area Search:

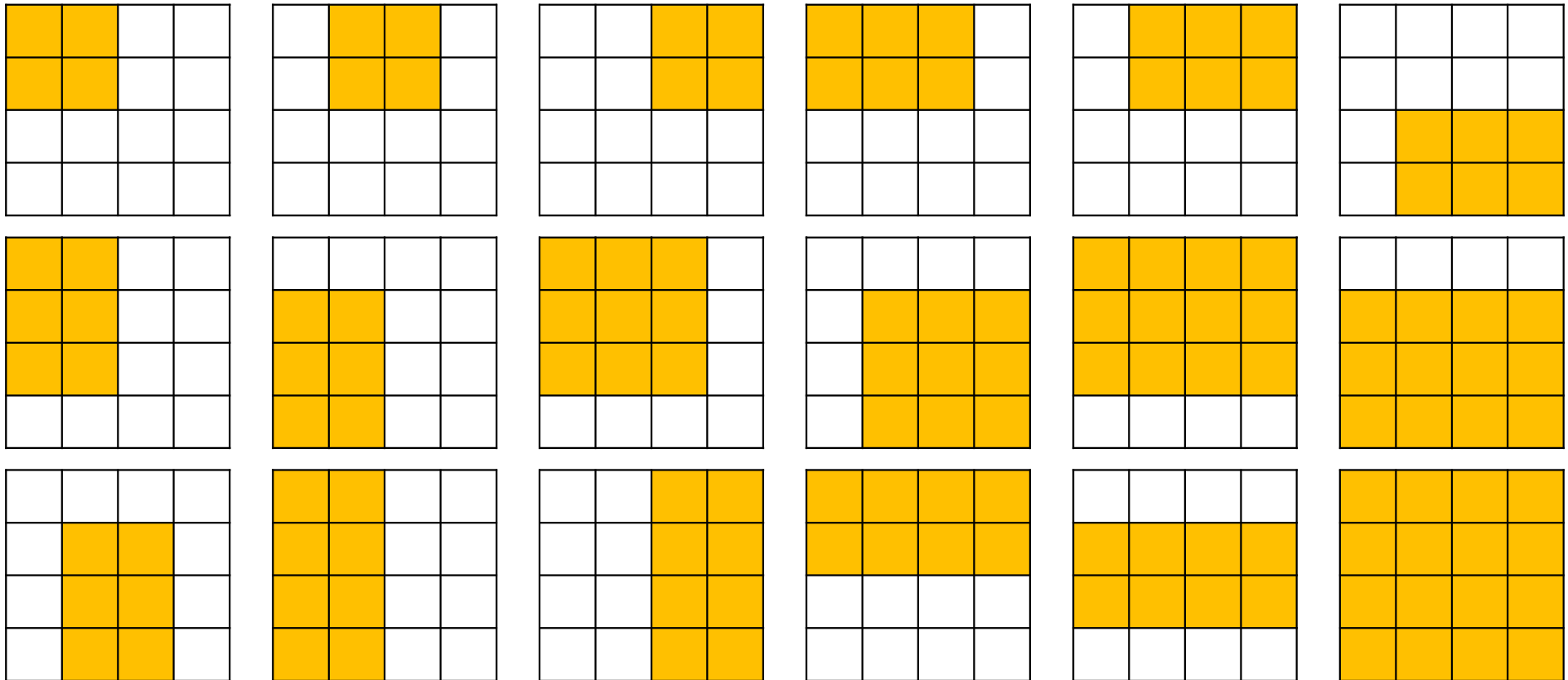
- Efficiency of similar region search

Assume:

$2 \leq \text{width} \leq 4$

$2 \leq \text{height} \leq 4$

$$O(n^2m^2)$$



02

区域相似性

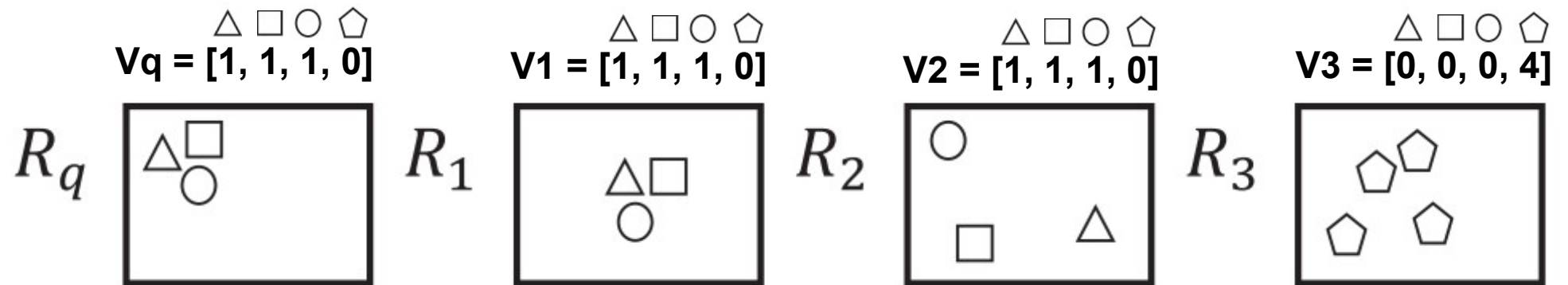


区域相似性



What affects the similarity:

- categories of POIs in the region



区域相似性

What affects the similarity:

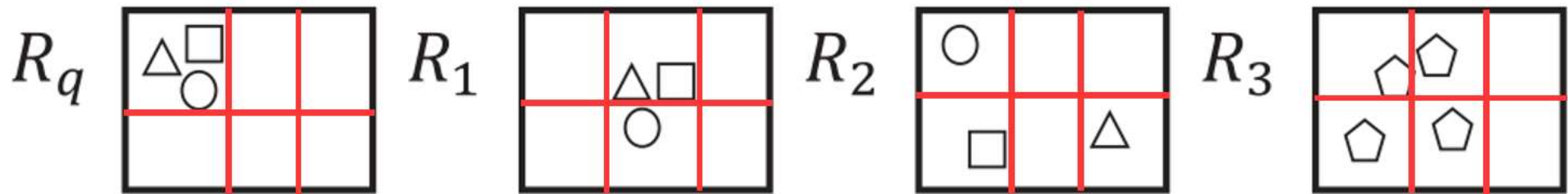
- Relative locations of the objects in a region

$$\begin{array}{c} \triangle \square \bigcirc \pentagon \\ Vq_{11} = [1, 1, 1, 0] \\ Vq_{12} = [0, 0, 0, 0] \\ Vq_{13} = [0, 0, 0, 0] \\ Vq_{21} = [0, 0, 0, 0] \\ Vq_{22} = [0, 0, 0, 0] \\ Vq_{23} = [0, 0, 0, 0] \end{array}$$

$$\begin{array}{c} \triangle \square \bigcirc \pentagon \\ V1_{11} = [0, 0, 0, 0] \\ V1_{12} = [1, 1, 0, 0] \\ V1_{13} = [0, 0, 0, 0] \\ V1_{21} = [0, 0, 0, 0] \\ V1_{22} = [0, 0, 1, 0] \\ V1_{23} = [0, 0, 0, 0] \end{array}$$

$$\begin{array}{c} \triangle \square \bigcirc \pentagon \\ V2_{11} = [0, 0, 1, 0] \\ V2_{12} = [0, 0, 0, 0] \\ V2_{13} = [0, 0, 0, 0] \\ V2_{21} = [0, 1, 0, 0] \\ V2_{22} = [0, 0, 0, 0] \\ V2_{23} = [1, 0, 0, 0] \end{array}$$

$$\begin{array}{c} \triangle \square \bigcirc \pentagon \\ V3_{11} = [0, 0, 0, 1] \\ V3_{12} = [0, 0, 0, 1] \\ V3_{13} = [0, 0, 0, 0] \\ V3_{21} = [0, 0, 0, 1] \\ V3_{22} = [0, 0, 0, 1] \\ V3_{23} = [0, 0, 0, 0] \end{array}$$



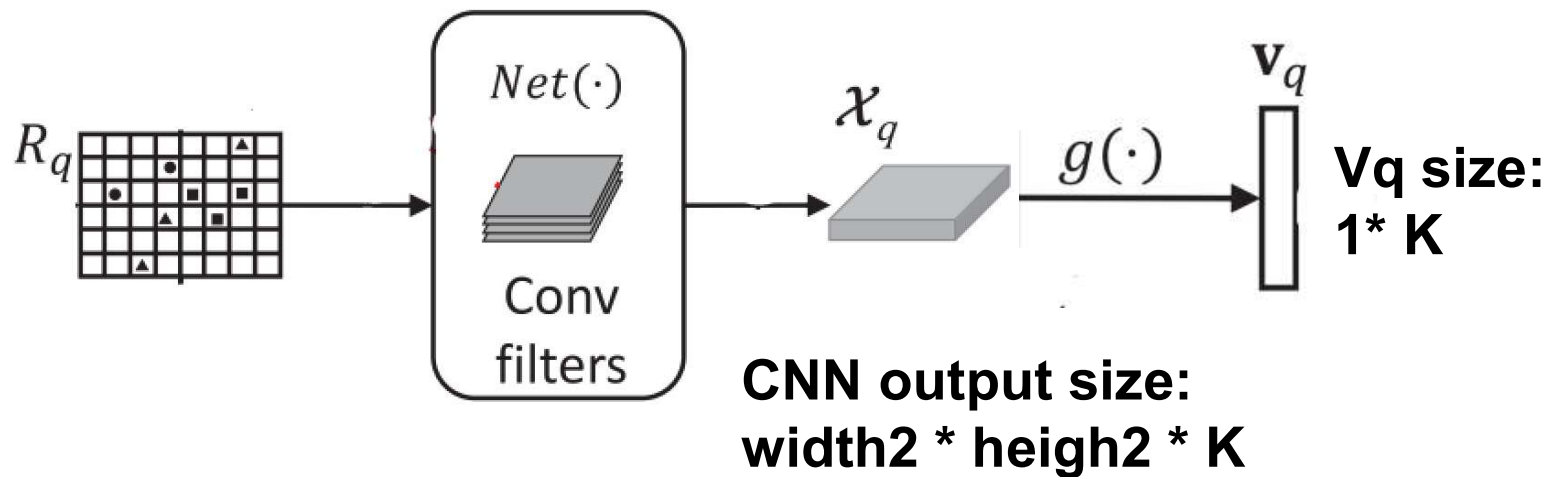
区域相似性

How to extract regional features:

- Relative locations of the objects in a region

Input size:

width1 * height1 * channel size

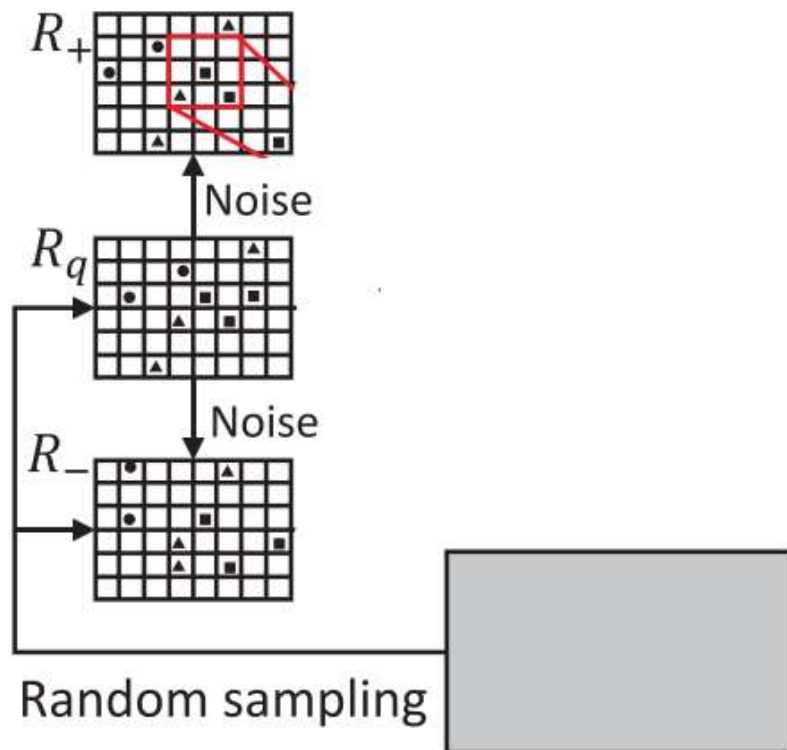


区域相似性



How to get the data whether the regions are similar or not:

- **Hand-made**

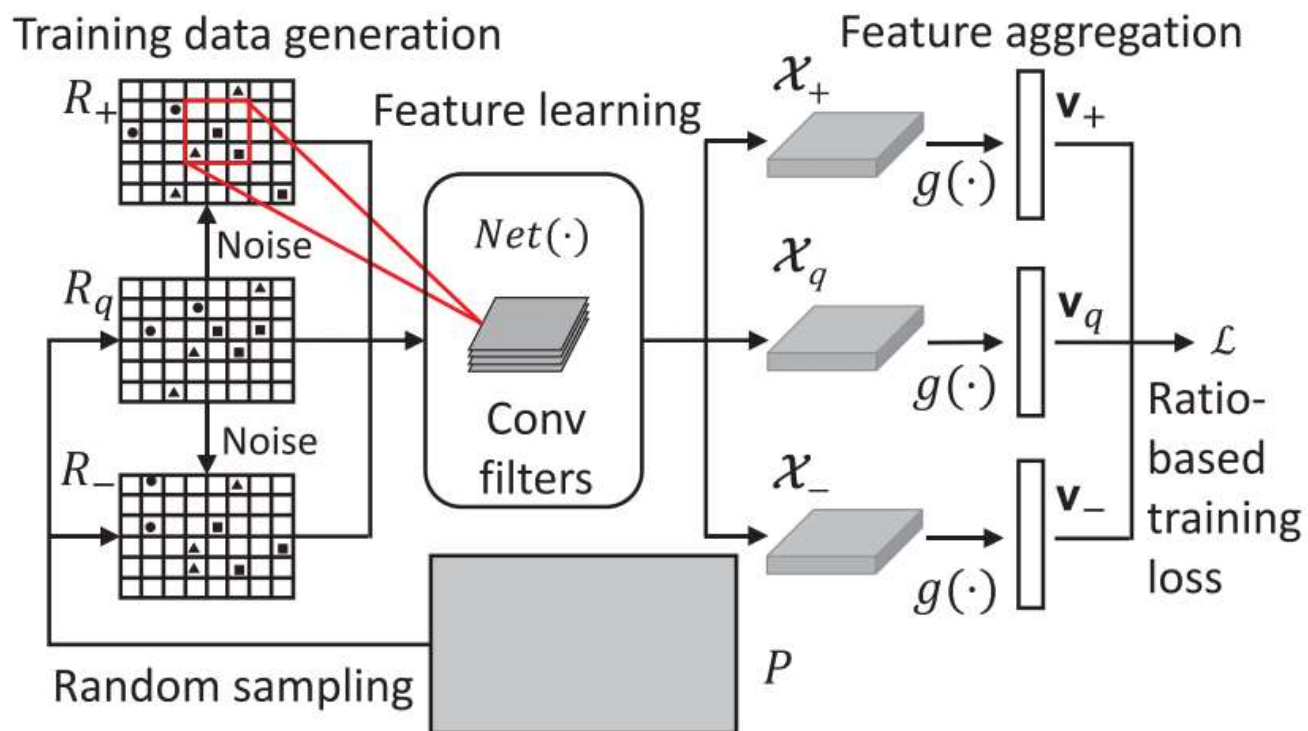


- randomly removing objects
- adding random objects at random locations
- randomly shifting the locations

区域相似性



Framework:



$g(\cdot)$: global max-pooling

$$\text{sim}(R_1, R_2) = \frac{1}{1 + \|\mathbf{v}_1 - \mathbf{v}_2\|_2}.$$

$$d_+ = \|\text{Net}(x_q) - \text{Net}(x_+)\|_2$$

$$d_- = \|\text{Net}(x_q) - \text{Net}(x_-)\|_2$$

$$\mathcal{L} = \sum_{(x_q, x_+, x_-)} \max\{0, d_+ - d_- + \delta\} + \lambda \|\text{Net}(\cdot)\|_2.$$

03

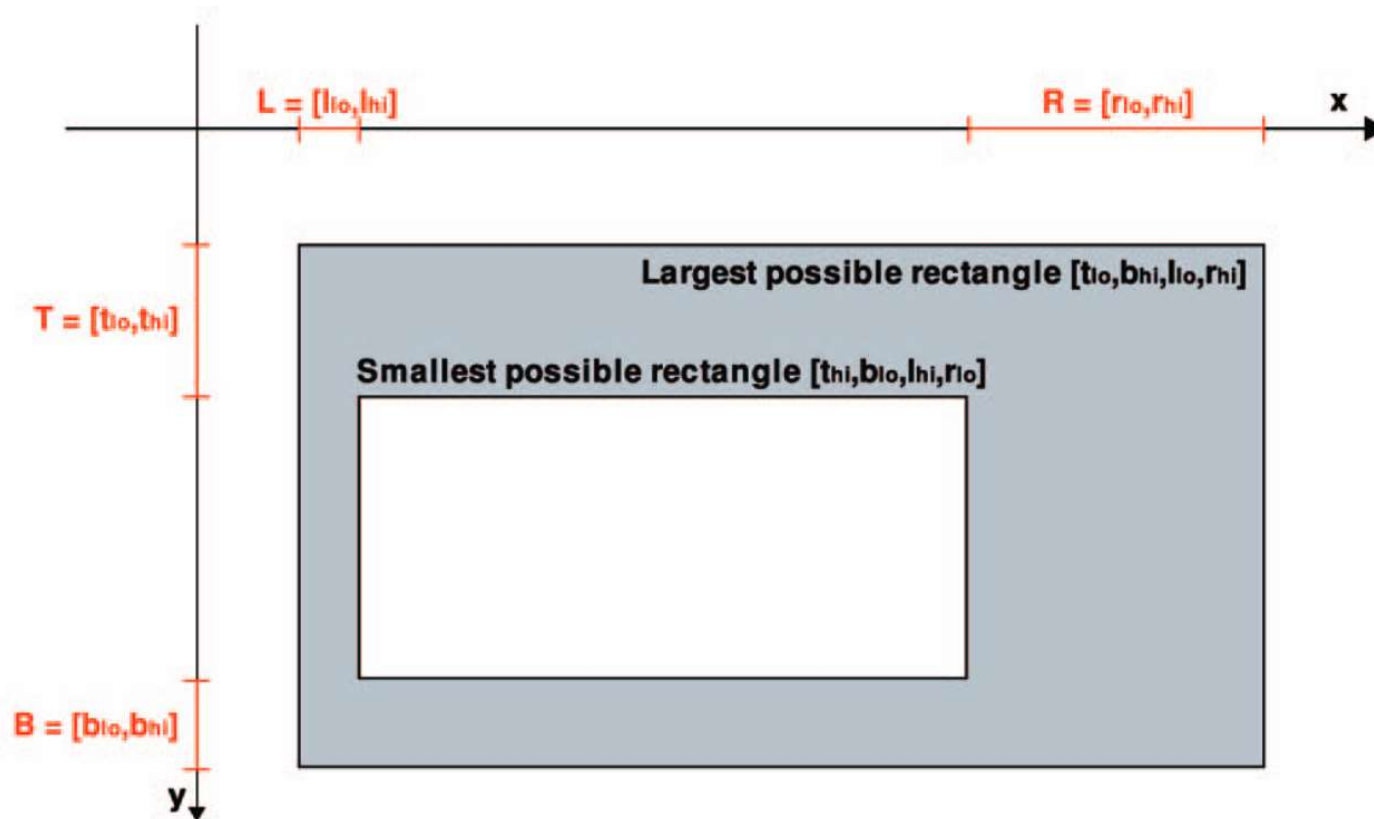
搜索算法



搜索算法



What are the candidate areas:

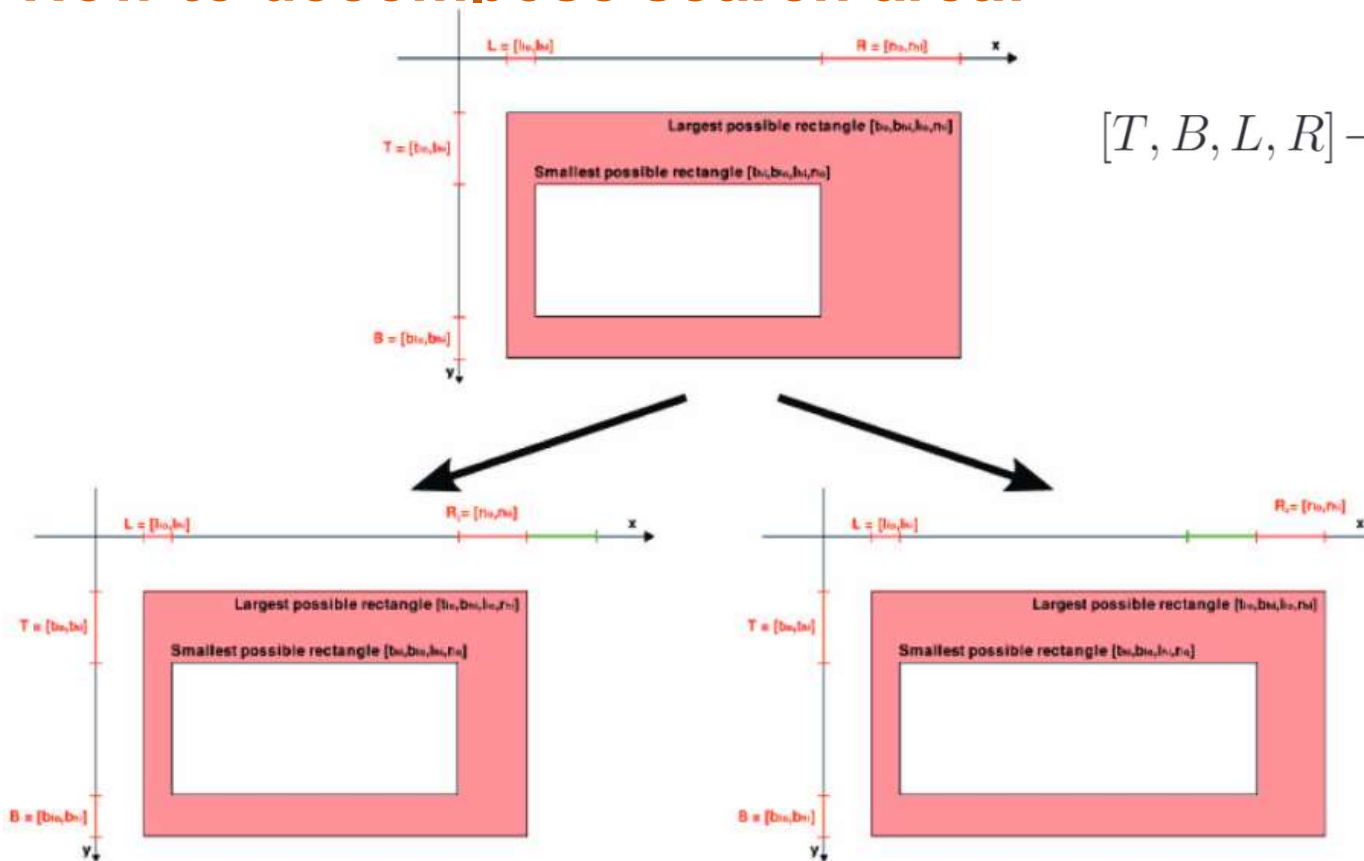


$$[T, B, L, R]$$

$$T = [t_{low}, t_{high}] \text{ etc.}$$

搜索算法

How to decompose search area:



$$[T, B, L, R] \rightarrow [T, B, L, R_1] \cup [T, B, L, R_2]$$

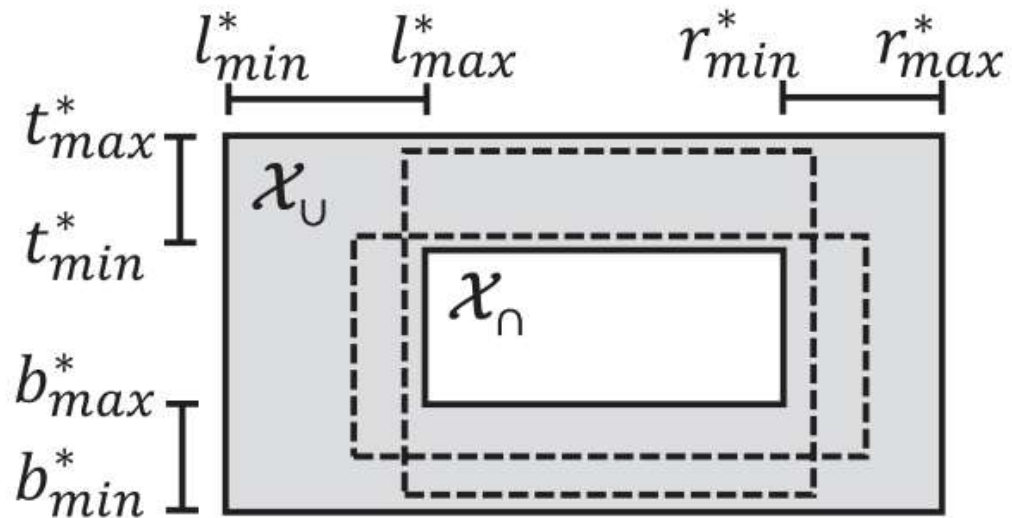
$$R_1 := [r_{lo}, \lfloor \frac{r_{lo} + r_{hi}}{2} \rfloor]$$

$$R_2 := [\lfloor \frac{r_{lo} + r_{hi}}{2} \rfloor + 1, r_{hi}]$$

搜索算法

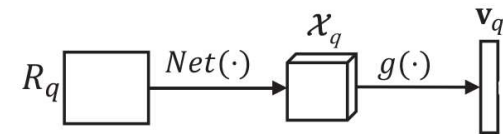


How to define the lower bound:



$$d = \sqrt{\sum_{k \in [1, K]} (v_q[k] - v_c[k])^2}$$

Recall: global pooling to represent region



global pooling is monotone increasing

$$X_n \subset X_c \subset X_U, \forall X_c \in S$$

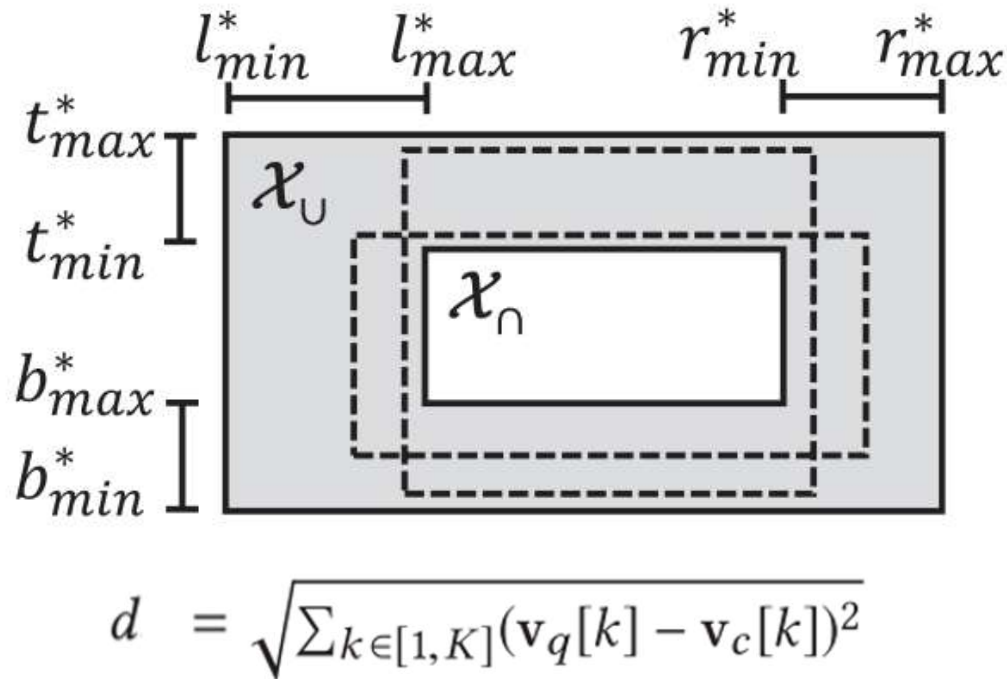
$$g(X_n) \leq g(X_c) \leq g(X_U)$$

$$\text{i.e., } v_n[k] \leq v_c[k] \leq v_U[k], \forall k \in [1, K]$$

搜索算法



How to define the lower bound:



Given: $\mathcal{X}_\cap \subset \mathcal{X}_c \subset \mathcal{X}_\cup, \forall \mathcal{X}_c \in \mathcal{S}$

$$v_\cap[k] \leq v_c[k] \leq v_\cup[k], \forall k \in [1, K]$$

It can be deduced that:

- (i) $(v_q[k] - v_c[k])^2 \geq (v_q[k] - v_\cap[k])^2$, if $v_\cap[k] \geq v_q[k]$.
- (ii) $(v_q[k] - v_c[k])^2 \geq (v_q[k] - v_\cup[k])^2$, if $v_\cup[k] \leq v_q[k]$.

Lower bound:

$$v_\cap[k] > v_q[k] \text{ as } k_1 \quad v_\cup[k] < v_q[k] \text{ as } k_2$$

$$\hat{f}(\mathcal{S}|\mathcal{X}_q) = \sqrt{\sum_{k_1} (v_q[k_1] - v_\cap[k_1])^2 + \sum_{k_2} (v_q[k_2] - v_\cup[k_2])^2}$$

搜索算法



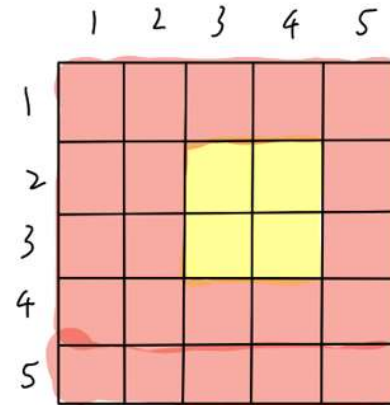
Search procedure:

Algorithm 1: ExactSFRS

Input: Initial search space S_V , query feature region X_q ,
distance lower bound $\hat{f}(\cdot)$

Output: top- N similar feature regions \mathcal{F}

```
1 begin
2    $\mathcal{F} \leftarrow \emptyset; Q \leftarrow \emptyset;$ 
3    $Q.Insert(S_V);$ 
4   repeat
5     repeat
6        $S' \leftarrow Q.RetrieveTop();$ 
7       split  $S \rightarrow S_1 \cup S_2;$ 
8        $Q.Insert((\hat{f}(S_1|X_q), S_1));$ 
9        $Q.Insert((\hat{f}(S_2|X_q), S_2));$ 
10    until  $|S'| = 1;$ 
11     $\mathcal{F} \leftarrow \mathcal{F} \cup S';$ 
12  until  $|\mathcal{F}| = N;$ 
13 end
```



$L [1,2]$
 $R [5,5]$
 $B [4,5]$
 $T [1,1]$
 $f = 4$

Queue

$L[1,2] \ R[5,5] \ B[4,5] \ T[1,1] \ f=4$

搜索算法



Search procedure:

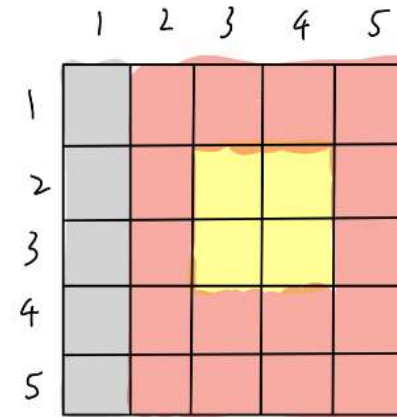
Algorithm 1: ExactSFRS

Input: Initial search space S_V , query feature region X_q , distance lower bound $\hat{f}(\cdot)$

Output: top- N similar feature regions \mathcal{F}

```

1 begin
2    $\mathcal{F} \leftarrow \emptyset; Q \leftarrow \emptyset;$ 
3    $Q.Insert(S_V);$ 
4   repeat
5     repeat
6        $S' \leftarrow Q.RetrieveTop();$ 
7       split  $S \rightarrow S_1 \cup S_2;$ 
8        $Q.Insert((\hat{f}(S_1|X_q), S_1));$ 
9        $Q.Insert((\hat{f}(S_2|X_q), S_2));$ 
10    until  $|S'| = 1;$ 
11     $\mathcal{F} \leftarrow \mathcal{F} \cup S';$ 
12  until  $|\mathcal{F}| = N;$ 
13 end
  
```



$L [2,2]$
 $R [5,5]$
 $B [4,5]$
 $T [1,1]$
 $f = 5$

Queue

$L [2,2] \quad R [5,5] \quad B [4,5] \quad T [1,1] \quad f = 5$
 $L [1,1] \quad R [5,5] \quad B [4,5] \quad T [1,1] \quad f = 6$

搜索算法



Search procedure:

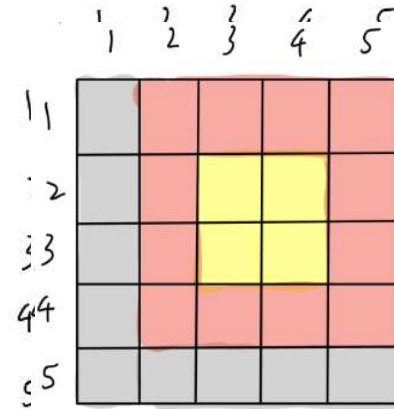
Algorithm 1: ExactSFRS

Input: Initial search space S_V , query feature region X_q , distance lower bound $\hat{f}(\cdot)$

Output: top- N similar feature regions \mathcal{F}

```

1 begin
2    $\mathcal{F} \leftarrow \emptyset; Q \leftarrow \emptyset;$ 
3    $Q.Insert(S_V);$ 
4   repeat
5     repeat
6        $S' \leftarrow Q.RetrieveTop();$ 
7       split  $S \rightarrow S_1 \cup S_2;$ 
8        $Q.Insert((\hat{f}(S_1|X_q), S_1));$ 
9        $Q.Insert((\hat{f}(S_2|X_q), S_2));$ 
10    until  $|S'| = 1;$ 
11     $\mathcal{F} \leftarrow \mathcal{F} \cup S';$ 
12  until  $|\mathcal{F}| = N;$ 
13 end
  
```



$L [2,2] ;$
 $R [5,5] ;$
 $B [4,4] ;$
 $T [1,1] ;$
 $f = 8$

Queue

$L [2,2] R [5,5] B [4,4] T [1,1] f = 5.5$
 $L [1,1] R [5,5] B [4,4] T [1,1] f = 6$
 $L [1,1] R [5,5] B [4,4] T [1,1] f = 8$

搜索算法



Search procedure:

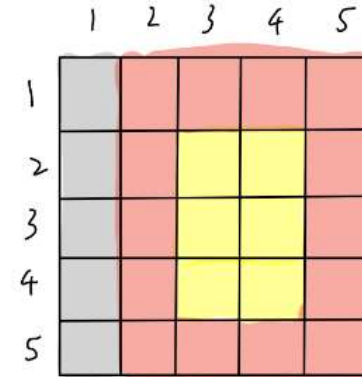
Algorithm 1: ExactSFRS

Input: Initial search space S_V , query feature region X_q , distance lower bound $\hat{f}(\cdot)$

Output: top- N similar feature regions \mathcal{F}

```

1 begin
2    $\mathcal{F} \leftarrow \emptyset$ ;  $Q \leftarrow \emptyset$ ;
3    $Q.Insert(S_V)$ ;
4   repeat
5     repeat
6        $S' \leftarrow Q.RetrieveTop()$ ;
7       split  $S \rightarrow S_1 \cup S_2$ ;
8        $Q.Insert((\hat{f}(S_1|X_q), S_1))$ ;
9        $Q.Insert((\hat{f}(S_2|X_q), S_2))$ ;
10    until  $|S'| = 1$ ;
11     $\mathcal{F} \leftarrow \mathcal{F} \cup S'$ ;
12  until  $|\mathcal{F}| = N$ ;
13 end
    
```



$L[2,2]$
 $R[5,5]$
 $B[5,5]$
 $T[1,1]$
 $f=5.5$

Queue

$L[2,2]$ $R[5,5]$ $B[5,5]$ $T[1,1]$ $f=5.5$
 $L[1,1]$ $R[5,5]$ $B[4,5]$ $T[1,1]$ $f=6$
 $L[1,2]$ $R[5,5]$ $B[4,4]$ $T[1,1]$ $f=8$



an area found!

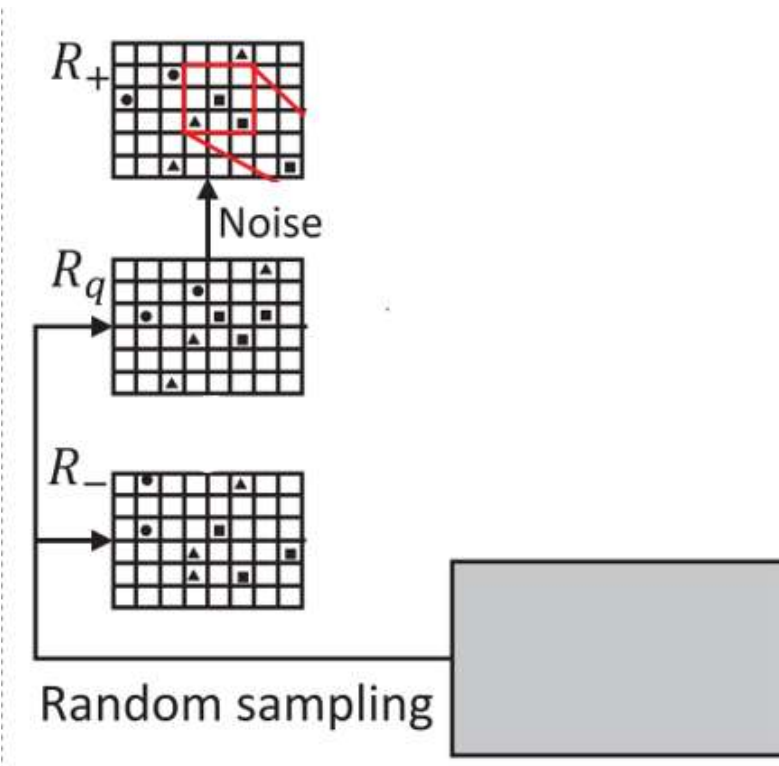
04

实验



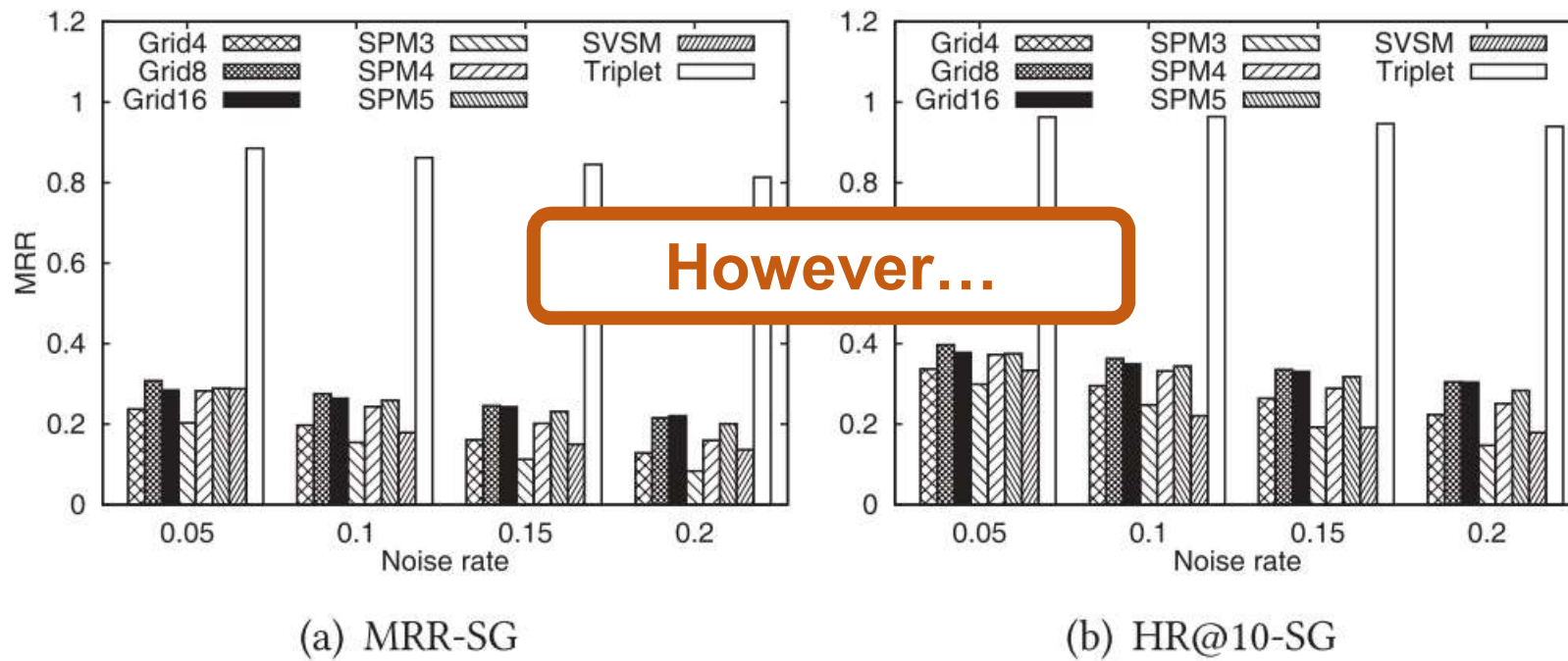
How to get the test data:

- Also hand-made



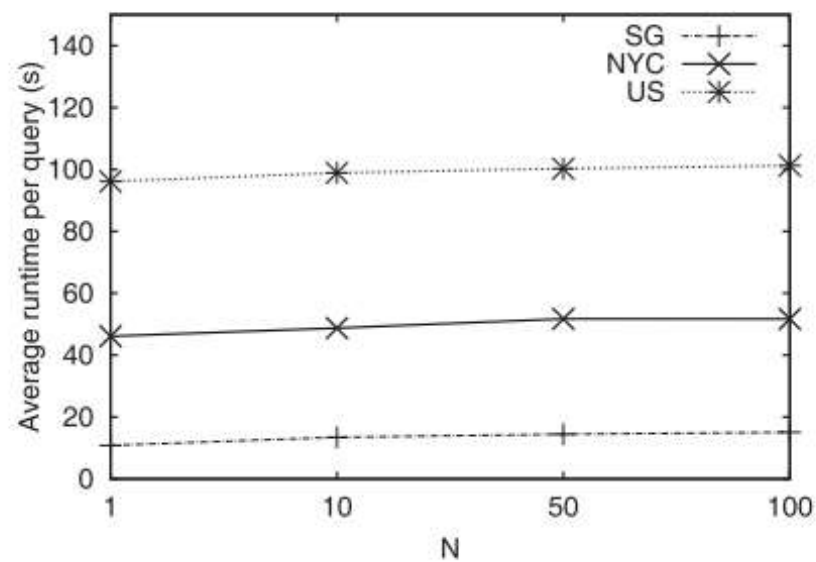
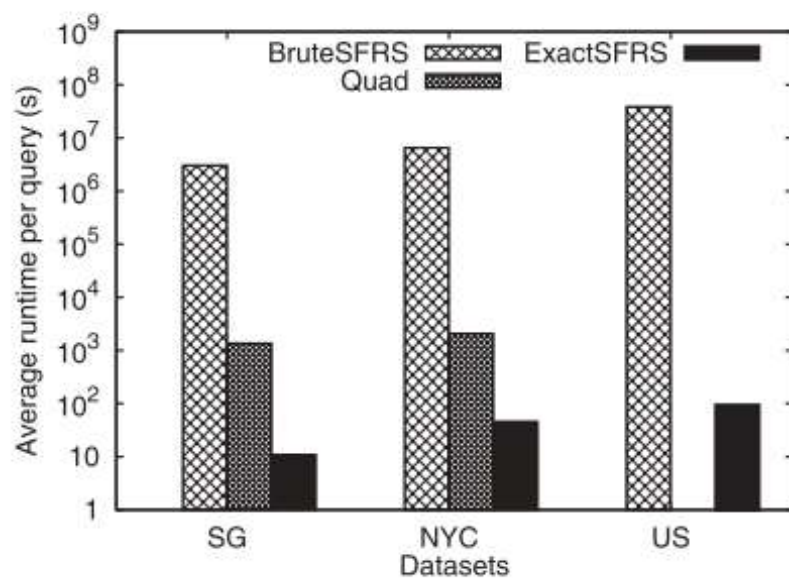
- randomly sample **2000 regions** that contains more than **50 objects** as the test queries
- The height and width of a region vary from **640m to 3km**
- candidate regions are constrained by **$0.5w_q \leq w_c \leq 2w_q$** and **$0.5h_q \leq h_c \leq 2h_q$**

Effectiveness:



实验

Efficiency:



05

讨论





还可以做点什么...区域相似性

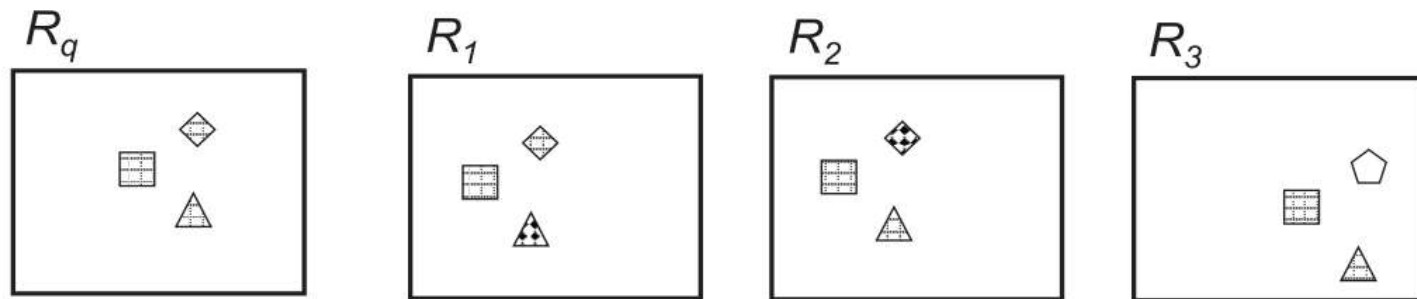
What else can be done:

- relations between categories[1]

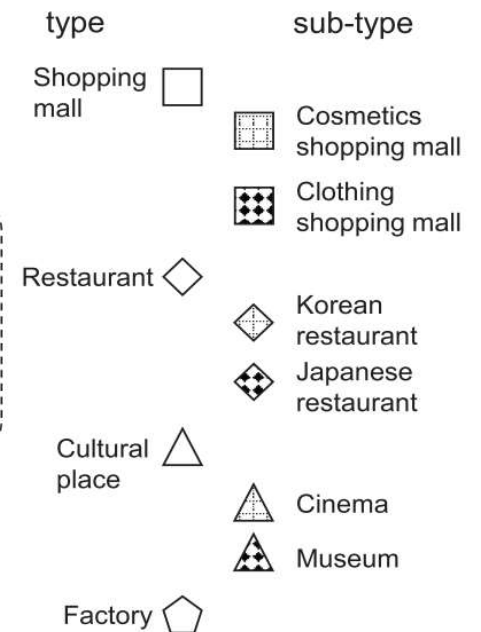
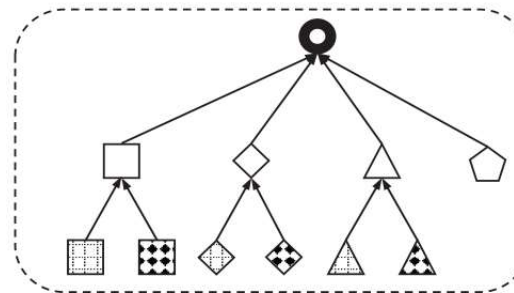
$$\text{sim}(R_q, R_1) = \text{sim}(R_q, R_2) = \text{sim}(R_q, R_3)$$

due to the one-hot embedding of types

It should be $\text{sim}(R_q, R_2) > \text{sim}(R_q, R_3)$



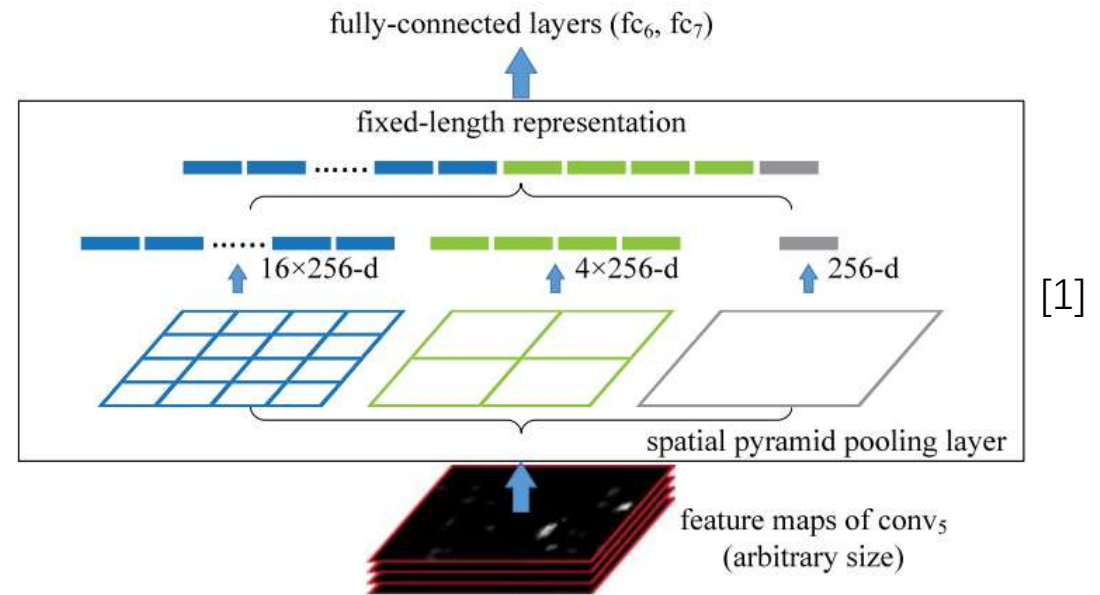
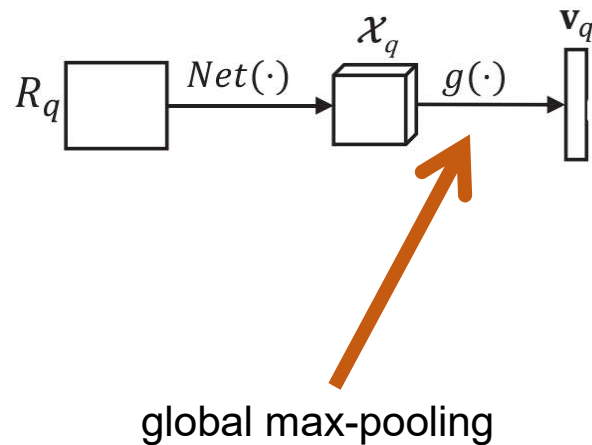
Type Hierarchy



还可以做点什么...区域相似性

What else can be done:

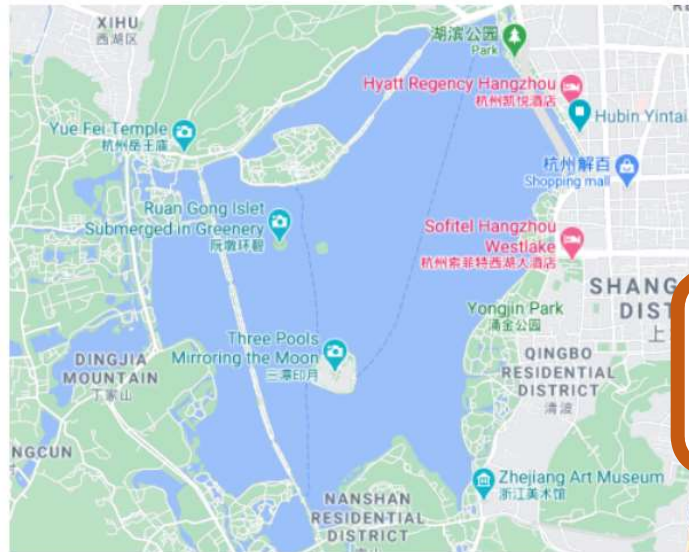
- feature aggregation method within a grid



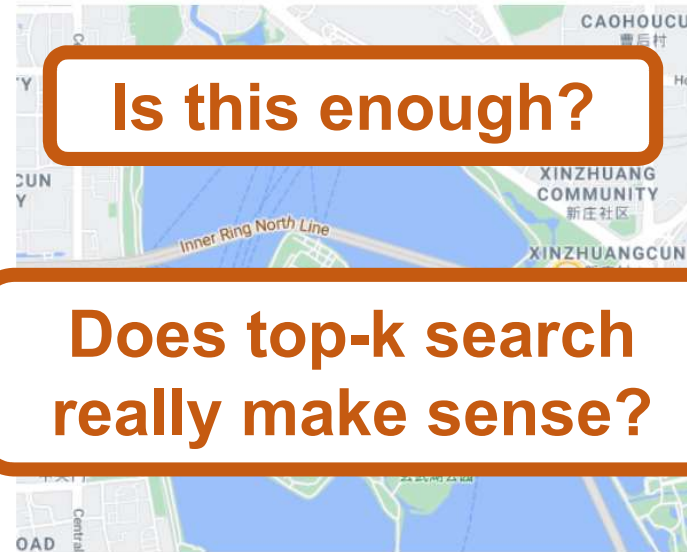
还可以做点什么...区域相似性

What else can be done:

- Dynamic surroundings impact --- What kind of area do users really want ?



(a) Hangzhou - West Lake



(b) Nanjing - Xuanwuhu Park



(c) Nanjing - XinJieKou

Is this enough?

Does top-k search really make sense?



谢谢大家

张斌杰 bj_zhang@seu.edu.cn



Gao Cong (丛高)

Block N4, 2c-103
School of Computer Science and Engineering
Nanyang Technological University
Email gaocong@ntu.edu.sg

1. Querying and Exploring Geospatial Data

1.1 Querying spatio-textual (geo-textual) data streams

Selected publications:

- SSTD: A Distributed System on Streaming Spatio-Textual Data, PVLDB 2020
- STAR: A Distributed Stream Warehouse System for Spatial Data. SIGMOD Conference 2020: 2761-2764 (Demo)
- Distributed Publish/Subscribe Query Processing on the Spatio-Textual Data Stream (ICDE 17)
- Diversity-aware top-k publish/subscribe on text stream (SIGMOD 15)
- Temporal spatial-keyword top-k publish/subscribe on geo-textual data stream (ICDE15 and demo in [VLDB14](#))
- Boolean spatial-keyword publish/subscribe on geo-textual data stream ([SIGMOD13](#) and demo in [VLDB14](#))

1.2 Data exploration for spatial data: Region search & topic exploration

Selected publications:

- SURGE: Continuous Detection of Bursty Regions Over a Stream of Spatial Objects (TKDE19, ICDE18)
- Finding attribute-aware similar regions for data analysis (PVLDB 19)
- Efficient Similar Region Search with Deep Metric Learning (KDD 18)
- Efficient Selection of Geospatial Data on maps for Interactive Visualized Exploration (SIGMOD 18)
- Towards Best Region Search for Data Exploration (SIGMOD 2016)
- Topic Exploration in Spatio-Temporal Document Collections(SIGMOD 2016, VLDBJ19)

1.2 Spatial keyword queries

- On Spatial Pattern Matching (ICDE'17, VLDBJ'19)
- Answering the m-closest keywords query (SIGMOD 15)
- Search regions of interest for user exploration ([VLDB14](#))
- Distributed spatial keyword querying on road networks ([EDBT14](#))
- An evaluation of 12 geo-spatial indexes ([VLDB13](#)). Code available [here](#).
- An overview paper on spatial-keyword querying ([invited paper in ER](#))
- Route planning: answering queries like “a most popular route such that it passes by shopping malls, restaurant, and p
- Efficient processing of several types of spatial keyword queries ([VLDB09](#), [PVLDB10](#), [SIGMOD11a](#)). Code for our SI
- TODS
- Efficient algorithms and cost models for reverse spatial-keyword k-nearest neighbor search ([SIGMOD11b](#), TODS14)
- Efficient spatial keyword search in trajectory databases ([unpublished paper](#))

Gao Cong (丛高)

Block N4, 2c-103
School of Computer Science and Engineering
Nanyang Technological University
Email gaocong@ntu.edu.sg

2. Spatial Data Mining and Spatial-temporal Data Mining

2.1 Intelligent transportation using trajectory data

Selected Publications:

- Online Anomalous Trajectory Detection with Deep Generative Sequence Modeling, ICDE2020
- Spatial Transition Learning on Road Networks with Deep Probabilistic Models, ICDE 2020
- Learning Travel Time Distributions with Deep Generative Model. (WWW 2019)

2.1 Data driven smart city applications

- Periodic-CRN: A Convolutional Recurrent Model for Crowd Density Prediction with Recurring Periodic Patterns (IJCAI, 2018)
- Efficient Similar Region Search with Deep Metric Learning (KDD 2018)

2.3 Spatial graph mining, POI recommendation & prediction

- Densely Connected User Community and Location Cluster Search in Location-Based Social Networks, SIGMOD2020
- Context-aware Deep Model for Joint Mobility and Time Prediction, WSDM 2020

4. Recommendation, POI recommendation and User Behaviour Modeling

4.1 Recommendation and group recommendation

Selected publications:

- HyperML: A Boosting Metric Learning Approach in Hyperbolic Space for Recommender Systems. WSDM 2020 (**Best paper award runner-up**)
- Global Context Enhanced Graph Neural Networks for Session-based Recommendation, SIGIR 2020
- Interact and Decide: Medley of Sub-Attention Networks for Effective Group Recommendation (SIGIR 19)
- Group Recommendation based on topic models(KDD14)

4.2 POI recommendation

- HME: A Hyperbolic Metric Embedding Approach for Next-POI Recommendation, SIGIR 2020
- A new POI recommendation approach, which performs better than previous approaches in experiments (SIGIR 2015)
- SAR: A sentiment-aspect-region model for user preference analysis and POI/user recommendation. The model provides explanations for recommendation results. (ICDE 2015)
- A general graph model for recommendation in heterogeneous networks and its applications in event-based social networks (ICDE 2015)
- Diversity-aware POI recommendation (AAAI 2015)
- Time-aware POI recommendation ([SIGIR13](#), CIKM14). Datasets available [here](#)
- Mining significant semantic locations from user generated GPS data for recommendation ([PVLDB10](#))

4.3 User behaviour modeling

- W4: Discovering spatio-temporal topics for individual users and its various applications, e.g., requirement-aware POI recommendation ([KDD13](#), TOIS15). Datasets available [here](#)