

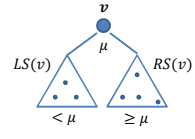
## VP-tree in Search of Nearest Neighbors within a Given Radius

Marzena Kryszkiewicz  
Instytut Informatyki  
Politechnika Warszawska

### The Idea of Constructing a VP-Tree

- A node in VP-Tree contains:**

- $v \in D$
- $\mu = \text{median}(\{u \in S(v) | \text{distance}(u, v)\})$ , where  $S(v)$  is the subtree rooted in  $v$
- link to left subtree  $LS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) < \mu\}$
- link to right subtree  $RS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) \geq \mu\}$



- Each point in  $D$  is stored only once in VP-Tree.**

- Idea of how to select a point from  $D$  to be stored in the root of the VP-Tree:**

- A point in the root of the VP-tree, say point  $v$ , should be the one with the maximal spread of its distances to all points in  $D$ .

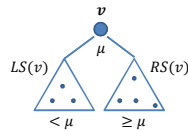
- Idea of how to select a point to be the root of a subtree covering a subset  $D'$  of points in  $D$ :**

- A point in the root of this subtree, say point  $v$ , should be the one with the maximal spread of its distances to all points in  $D'$ .

### Practical Construction of a VP-Tree

- A node in VP-Tree contains:**

- $v \in D$
- $\mu = \text{median}(\{u \in S(v) | \text{distance}(u, v)\})$ , where  $S(v)$  is the subtree rooted in  $v$
- link to left subtree  $LS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) < \mu\}$
- link to right subtree  $RS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) \geq \mu\}$



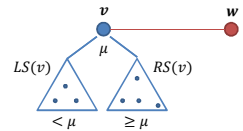
- Practical selection of a point from (a subset  $D'$  of)  $D$  to be stored in the root of a (sub-)tree:**

- A random sample  $S_1$  of points from (subset  $D'$  of)  $D$  constitutes a set of candidates to be stored in the root of the (sub-)tree.
- Their spreads of distances are calculated with respect to another random sample  $S_2$  of points from (subset  $D'$  of)  $D$ .
- The candidate point  $v$  from  $S_1$  with the maximal spread of its distances to the points in the sample  $S_2$  is stored in the root of the (sub-)tree.
- The real median of this point  $v$  is calculated based on its distances to all points in (subset  $D'$  of)  $D$ , and is also stored in the root of the (sub-)tree.

### $k/k^+$ -NN Search in VP-Tree

- A node in VP-Tree contains:**

- $v \in D$
- $\mu = \text{median}(\{u \in S(v) | \text{distance}(u, v)\})$ , where  $S(v)$  is the subtree rooted in  $v$
- link to left subtree  $LS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) < \mu\}$
- link to right subtree  $RS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) \geq \mu\}$



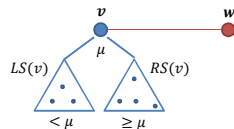
- Search for  $k/k^+$ -NN of point  $u$  within  $\epsilon$  radius in node  $v$  of VP-Tree:**

- $\text{distance}(w, v)$ ,
- Cond. 1:  $\text{distance}(w, v) - \mu \geq \epsilon$ . If true, then for each point  $u$  in  $LS(v)$ ,  $\text{distance}(w, v) - \text{distance}(u, v) > \text{distance}(w, v) - \mu \geq \epsilon$ . Thus,  $\text{distance}(w, v) - \text{distance}(u, v) > \epsilon$ , so  $LS(v)$  does not contain  $k/k^+$ -NN( $w$ ) within the  $\epsilon$  radius.
- Cond. 2:  $\mu - \text{distance}(w, v) > \epsilon$ . If true, then for each point  $u$  in  $RS(v)$ ,  $\text{distance}(u, v) - \text{distance}(w, v) \geq \mu - \text{distance}(w, v) > \epsilon$ . Thus,  $\text{distance}(u, v) - \text{distance}(w, v) > \epsilon$ , so  $RS(v)$  does not contain  $k/k^+$ -NN( $w$ ) within the  $\epsilon$  radius.

### Improved $k/k^+$ -NN Search in VP-Tree...

- A node in VP-Tree contains:**

- $v \in D$
- $\mu = \text{median}(\{u \in S(v) | \text{distance}(u, v)\})$ , where  $S(v)$  is the subtree rooted in  $v$
- link to left subtree  $LS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) < \mu\}$
- link to right subtree  $RS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) \geq \mu\}$



- Search for  $k/k^+$ -NN of point  $u$  within  $\epsilon$  radius in node  $v$  of VP-Tree:**

- $\text{distance}(w, v)$ ,
- Cond. 1:  $\text{distance}(w, v) - \mu \geq \epsilon$ . If true, then for each point  $u$  in  $LS(v)$ ,  $\text{distance}(w, v) - \text{distance}(u, v) > \epsilon$ , so  $kNN(w)$  is not in  $LS(v)$  within  $\epsilon$  radius.
- Cond. 2:  $\mu - \text{distance}(w, v) > \epsilon$ . If true, then for each point  $u$  in  $RS(v)$ ,  $\text{distance}(u, v) - \text{distance}(w, v) > \epsilon$ , so  $kNN(w)$  is not in  $RS(v)$  within  $\epsilon$  radius.

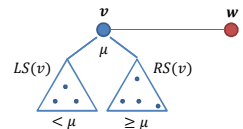
- Improved search for  $k/k^+$ -NN of point  $u$  within  $\epsilon$  radius in node  $v$  of VP-Tree:**

- $\text{distance}(w, v)$ ,
- Cond. 1':  $\text{distance}(w, v) - \text{left\_bound} > \epsilon$ , where  $\text{left\_bound}$  is the maximum of the distances from point  $v$  to all points in  $LS(v)$ .
- Cond. 2':  $\text{right\_bound} - \text{distance}(w, v) > \epsilon$ , where  $\text{right\_bound}$  is the minimum of the distances from point  $v$  to all points in  $RS(v)$ .

### Improved $k/k^+$ -NN Search in VP-Tree

- A node in VP-Tree contains:**

- $v \in D$
- $\mu = \text{median}(\{u \in S(v) | \text{distance}(u, v)\})$ , where  $S(v)$  is the subtree rooted in  $v$
- link to left subtree  $LS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) < \mu\}$
- link to right subtree  $RS(v)$  embracing  $\{u \in S(v) \setminus \{v\} | \text{distance}(u, v) \geq \mu\}$



- Improved search for  $k/k^+$ -NN of point  $u$  within  $\epsilon$  radius in node  $v$  of VP-Tree:**

- $\text{distance}(w, v)$ ,
- Cond. 1':  $\text{distance}(w, v) - \text{left\_bound} > \epsilon$ , where  $\text{left\_bound}$  is the maximum of the distances from point  $v$  to all points in  $LS(v)$ .
- Cond. 2':  $\text{right\_bound} - \text{distance}(w, v) > \epsilon$ , where  $\text{right\_bound}$  is the minimum of the distances from point  $v$  to all points in  $RS(v)$ .

- Example.** Let  $\epsilon = 1$ ,  $\mu = 10.5$ ,  $\text{left\_bound}(v) = 8.5$ ,  $\text{right\_bound}(v) = 12$  and  $\text{distance}(w, v) = 10$ . Then,  $\text{distance}(w, v) - \text{left\_bound}(v) > \epsilon$  and  $\text{right\_bound}(v) - \text{distance}(w, v) > \epsilon$ , which means that neither  $LS(v)$  nor  $RS(v)$  contains any nearest neighbor of  $w$  within the  $\epsilon$  radius.

## References

- Kryszkiewicz M., Janczak B.: Basic Triangle Inequality Approach Versus Metric VP-Tree and Projection in Determining Euclidean and Cosine Neighbors. *Intelligent Tools for Building a Scientific Information Platform 2014*: 27-49
- Moore, A. W.: The Anchors Hierarchy: Using the Triangle Inequality to Survive High Dimensional Data. In: *Proc. of UAI, Stanford (2000)* 397–405
- Yanilos P. N.: Data Structures and Algorithms of Nearest Neighbor Search in General Metric Spaces. *Materiały z 4th ACM-SIAM Symposium on Discrete Algorithms*, 1993, 311-321
- Zezula, P., Amato, G., Dohnal, V., Bratko, M.: *Similarity Search: The Metric Space Approach*. Springer (2006)