

# Geometric Aspects and Auxiliary Features to Top-k Processing

## [Advanced Seminar]

Kyriakos Mouratidis

School of Information Systems  
Singapore Management University  
80 Stamford Road, Singapore 178902  
kyriakos@smu.edu.sg

**Abstract**—Top- $k$  processing is a well-studied problem with numerous applications that is becoming increasingly relevant with the growing availability of recommendation systems and decision making software on PCs, PDAs and smart-phones. The objective of this seminar is twofold. First, we will delve into the geometric aspects of top- $k$  processing. Second, we will cover complementary features to top- $k$  queries that have a strong geometric nature. The seminar will close with insights in the effect of dimensionality on the meaningfulness of top- $k$  queries, and interesting similarities to nearest neighbor search.

## 1. Introduction

Consider a dataset that contains a large number of *options* (e.g., restaurants, hotels, etc.). Each option  $\mathbf{r}$  has  $d$  *attributes*. In an example where the dataset contains hotels, the attributes could correspond to the ratings of the hotels on  $d$  different aspects, such as service, sleep quality, convenience of location, etc. The top- $k$  query is a common means to shortlist the  $k$  best options according to the user's preferences on the  $d$  data attributes. Specifically, in the most prevalent top- $k$  model, the user specifies a *query vector*  $\mathbf{q}$  which comprises a numeric weight for each attribute [7]. The score of an option is defined as the weighted sum of its attributes (equivalently, the dot product  $\mathbf{r} \cdot \mathbf{q}$ ), which in turn imposes a ranking among the available options. The  $k$  highest ranking options form the top- $k$  result and are reported to the user.

Despite its algebraic definition, top- $k$  processing has a strong geometric nature and a connection to traditional computational geometry problems. For example, if the options are treated as points in a  $d$ -dimensional space, top- $k$  computation can be seen as a sweeping of the data space from its top corner to the origin with a hyper-plane (normal to the query vector  $\mathbf{q}$ ) until  $k$  options are swept. In addition to ideas for query processing, this parallel reveals important properties of the problem, such as the fact that the top option in a dataset for any query vector lies on its *convex hull*.

Things become more interesting when variants or auxiliary features to top- $k$  processing are considered in the query space, i.e., the space where the query vector may

lie. Geometric properties in that space, and particularly the concept of *k-levels* from computational geometry, can be used for the efficient processing of ad-hoc top- $k$  queries over data streams [5] or the processing of *continuous* top- $k$  queries [14].

Furthermore, insights in the properties of the query space have given rise to very useful, complementary features (and measures) relevant to top- $k$  processing. An example is the association of the top- $k$  result with a region around the query vector  $\mathbf{q}$  (in query space) where the result remains the same [11], [15]. The volume of that region can be used as a measure for result sensitivity, while the region itself as a means for computation sharing among different top- $k$  queries (result caching), for exploratory analysis, etc. Another example is the computation of the maximum possible rank that an option could achieve, given the competition (i.e., the alternative options in the data set) [9]. This also entails the calculation of the exact regions of the query space where the maximum rank is achieved, which can be used for market impact analysis and customer profiling. The problem is related to *hyper-plane arrangements*, a very powerful concept in computational geometry [1], [2].

The first objective of this seminar is to shed some light on the connection between top- $k$  processing and standard computational geometry problems. In particular, we will review the key geometric concepts of (i) convex hull, (ii) half-space range reporting, (iii) hyperplane arrangement, and (iv)  $k$ -level, and we will explain how they can help us support top- $k$  queries, variants and auxiliary features.

The seminar will continue with a coverage of systems that exploit the aforementioned geometric fundamentals to efficiently process top- $k$  and related queries. Over 10 papers will be discussed, with a focus on methods that (i) are not limited to 2 dimensions only, and (ii) offer exact solutions. Regarding item (i), as we will explain, 2-dimensional solutions address degenerate versions of preference-based ranking, which may simplify processing but they also sacrifice the generality of the key ideas. Regarding (ii), approximate methods may be plausible for some problems, but they largely dismiss the geometric aspects of the problems in order to simplify them.

The seminar will conclude with insights into the effect of dimensionality on the behaviour and the usefulness

of top- $k$  ranking/processing. Connections to the standard *nearest neighbor query* (NN) will be drawn, and surprising similarities to its behaviour with dimensionality will be demonstrated.

The seminar is designed for 90 minutes. The presenter has not given a tutorial on this subject to another venue; the seminar is designed specifically with the MDM audience in mind. Due to its close connection to computational geometry, the seminar will be very “visual”, with a multitude of drawings used for illustration; this will hopefully make it eye-catching and fun (in addition to useful).

## 2. Target Audience and Relevance to MDM’16

The topic is very closely related to *spatial databases*, *multi-criteria decision making*, and *decision analytics*. The spatial database community, which forms a large fraction of the MDM audience, has a special interest in applied geometry and geometric reasoning at large scale, which is the core of this seminar. Given that the application domain of the seminar centers on preference-based querying and top- $k$  processing, the topic is directly related to audiences interested in recommendation systems, multi-criteria decision making and decision analytics.

Specific areas from the MDM’16 call for papers that match the seminar’s topic:

*Recommendations for Mobile Users:* The seminar centers on preference-based querying, and therefore lies at the core of recommendation systems. It follows a general formulation where preferences are expressed over a number of attributes; the latter could easily be location-dependent, such as the distances of the options from specific landmarks or from the user’s location.

*Data stream processing in mobile/sensor network:* Several of the papers reviewed, and their targeted applications, involve data streams and highly dynamic environments (e.g., [5], [14]). The sources of the data streams could be, among others, mobile or sensor networks.

*Indexing, Optimisation and Query Processing for moving objects/users:* The seminar will cover database papers, where the goal is scalability. Each and every system we will review exploits geometric properties with the very objective, exactly, to organise data with appropriate indices and to enable efficient query processing by smart algorithms and optimizations.

## 3. Seminar Outline

The seminar comprises 3 parts. The first covers standard computational geometry concepts and their relevance to top- $k$  processing, namely, convex hull, half-space range reporting, hyperplane arrangement, and  $k$ -level. The second, and the main part of the seminar, covers existing database work that exploits the connection of these problems to top- $k$  processing. The third part concludes the seminar by insights about the effect of dimensionality on the usefulness of top- $k$  and related queries, and draws a parallel to the traditional

NN query from spatial databases. Below we summarize the main systems to be reviewed in the second (and major) part of the seminar.

The first work covered is [5], where Das et al. consider the evaluation of ad-hoc top- $k$  queries over a data stream in the sliding window model. The main idea is that only a small subset of the options in the sliding window could appear in the top- $k$  result w.r.t. any query vector. To identify (and maintain) this small subset they rely on a geometric representation of the top- $k$  query and a notion of duality, where options and queries are mapped into lines and rays, respectively. Although tailored to 2 dimensions, this work involves fundamental principles that are key to the seminar, especially those regarding geometric arrangements.

Yu et al. [14] extend the principles of [5] to higher dimensions, targeting this time *continuous* top- $k$  queries. At the core of their approach lies the effective maintenance of the *query response surface*, which encodes the score and identity of the  $k$ -th result option for any query vector, and is therefore very relevant to  $k$ -levels.

Vlachou et al. [12], [13] study *reverse top- $k$  queries*. Even though their main focus is on the *bichromatic* version of the query (whose definition and treatment deviates from the geometric focus of this seminar), they introduce the *monochromatic* version too, and propose a geometric solution for 2 dimensions. Specified a focal option  $\mathbf{p}$  in a set of alternatives, they compute the parts of the query space where the query vector should lie so that  $\mathbf{p}$  belongs to the top- $k$  result. While applicable to 2 dimensions only, the proposed solution is an interesting use of  $k$ -levels for a top- $k$ -related problem.

Soliman et al. [11] consider uncertain scoring functions and identify the most “representative” top- $k$  result, under different definitions. First, they compute the most likely top- $k$  result if the query vector is randomly chosen. Next, they compute the top- $k$  result that is least dissimilar to all possible alternative results. Finally, they introduce sensitivity measures for a given top- $k$  result. Their approach relies on geometric insights and on operations involving hyperplane representations in query space.

Zhang et al. [15] introduce the concept of the *global immutable region* (GIR). The GIR is the maximal locus around a query vector  $\mathbf{q}$  where the top- $k$  result remains the same. The GIR is shown to be a convex polytope (in query space), produced by half-space intersection. To offer scalability, the authors rely on properties of the convex hull (and its connection to top- $k$  processing), while their most efficient algorithm borrows ideas from Clarkson’s algorithm, one of the most common methods for convex hull computation.

Mouratidis et al. [9] propose the *maximum rank query* (*MaxRank*). Given a focal option in a set of alternatives, *MaxRank* computes the highest rank this option may achieve under any possible user preference. It additionally reports all the regions in the query space where that rank is achieved. Its applications include market impact analysis, customer profiling, targeted advertising, etc. The proposed solution relies on hyperplane arrangements and a coverage counting problem therein.

He and Lo [6] consider *why-not* queries in the context of top- $k$  processing. Given a query vector  $\mathbf{q}$  and an option  $\mathbf{p}$  that does not belong to the top- $k$  result, the *why-not* query computes the smallest perturbation required in the query vector and/or in value  $k$  so that  $\mathbf{p}$  is included in the result. Even though the proposed solution relies on sampling (and is therefore approximate), it involves interesting geometric insights and relies on hyperplane arrangements in query space.

In addition to the above papers, we will make mention to studies of practical relevance, which however deviate from our focus, i.e., offer approximate answers, are bound to 2 dimensions, or have no geometric solutions. Examples include bichromatic reverse top- $k$  [13], reverse  $k$ -ranks [16] (which can be seen as the bichromatic version of *MaxRank*), external  $k$ -level computation in 2 dimensions [3],  $k$ -regret minimizing sets [4],  $k$ -hit [10], etc.

#### 4. Biography of the Presenter



Kyriakos Mouratidis was born in Greece in 1980. He completed his B.Sc. in Computer Science at Aristotle University of Thessaloniki (AUTH) in 2002, and his Ph.D. in Computer Science and Engineering at Hong Kong University of Science and Technology (HKUST) in 2006. In the same year, he joined the School of Information Systems at Singapore Management University (SMU), where he is currently an Associate Professor. His main research area is spatial databases, with a focus on continuous query processing, road network databases, and spatial optimization problems. His work in the last 3 years has concentrated on complementary features to top- $k$  queries, like for example [8], [15], [9]. A complete CV and publication list can be found at: <http://www.mysmu.edu/faculty/kyriakos/>

#### References

- [1] P. K. Agarwal and M. Sharir. Arrangements and their applications. In *Handbook of Computational Geometry*, pages 49–119. Elsevier, 1998.
- [2] M. D. Berg, O. Cheong, M. V. Kreveld, and M. Overmars. *Computational geometry: algorithms and applications*. Springer, 2008.
- [3] M. A. Cheema, Z. Shen, X. Lin, and W. Zhang. A unified framework for efficiently processing ranking related queries. In *EDBT*, pages 427–438, 2014.
- [4] S. Chester, A. Thomo, S. Venkatesh, and S. Whitesides. Computing  $k$ -regret minimizing sets. *PVLDB*, 7(5):389–400, 2014.
- [5] G. Das, D. Gunopulos, N. Koudas, and N. Sarkas. Ad-hoc top- $k$  query answering for data streams. In *VLDB*, pages 183–194, 2007.
- [6] Z. He and E. Lo. Answering *why-not* questions on top- $k$  queries. *IEEE Trans. Knowl. Data Eng.*, 26(6):1300–1315, 2014.
- [7] I. F. Ilyas, G. Beskales, and M. A. Soliman. A survey of top- $k$  query processing techniques in relational database systems. *ACM Comp. Surveys*, 40(4), 2008.
- [8] K. Mouratidis and H. Pang. Computing immutable regions for subspace top- $k$  queries. In *PVLDB*, pages 73–84, 2013.
- [9] K. Mouratidis, J. Zhang, and H. Pang. Maximum rank query. *PVLDB*, 8(12):1554–1565, 2015.
- [10] P. Peng and R. C. Wong.  $k$ -hit query: Top- $k$  query with probabilistic utility function. In *SIGMOD*, pages 577–592, 2015.
- [11] M. A. Soliman, I. F. Ilyas, D. Martinenghi, and M. Tagliasacchi. Ranking with uncertain scoring functions: semantics and sensitivity measures. In *SIGMOD*, pages 805–816, 2011.
- [12] A. Vlachou, C. Doulkeridis, Y. Kotidis, and K. Nørnvåg. Reverse top- $k$  queries. In *ICDE*, pages 365–376, 2010.
- [13] A. Vlachou, C. Doulkeridis, Y. Kotidis, and K. Nørnvåg. Monochromatic and bichromatic reverse top- $k$  queries. *IEEE Trans. Knowl. Data Eng.*, 23(8):1215–1229, 2011.
- [14] A. Yu, P. K. Agarwal, and J. Yang. Processing a large number of continuous preference top- $k$  queries. In *SIGMOD*, pages 397–408, 2012.
- [15] J. Zhang, K. Mouratidis, and H. Pang. Global immutable region computation. In *SIGMOD*, pages 1151–1162, 2014.
- [16] Z. Zhang, C. Jin, and Q. Kang. Reverse  $k$ -ranks query. *PVLDB*, 7(10):785–796, 2014.