# Two-way Recommendation System for Supervisor Selection using Historical Data and Skyband-view Queries

Annisa Annisa[1,*], Muhammad Rayhan Adyatma[1], Global Ilham Sampurno[1], Chen Li[2]

1  Department of Computer Science, Faculty of Mathematics and Natural Sciences, IPB University, Bogor, Indonesia
2  Graduate School of Informatics, Nagoya University, Chikusa, Nagoya 464-8602, Japan

| ARTICLE INFO | ABSTRACT |
|---|---|
| <br><br><br><br><br><br><br><br><br><br><br> | Completing the final project on time is one indicator of success in university. Unfortunately, many students have not been able to complete their studies on time. One of the factors that influence this is the failure in choosing the right supervisor. This study aims to create a recommendation system that can help students choose suitable supervisors. Unlike other studies, this research builds a two-way recommendation system that takes into account the preferences of students and supervisors. The previous study used the skyline views query concept to recommend dominant objects. However, the skyline view query concept has a major limitation: only skyline objects will be recommended to the user. Thus, students who are not skyline objects may not get a supervisor's recommendation, and vice versa. In this research, we use the concept of a skyband view query to overcome the limitation of the skyline view query. In addition to answering eight important queries from both parties, students and supervisors, the skyband view query concept is also able to overcome the shortcomings in previous research. Historical data from alumni is used to construct students' and supervisors' interests. This research succeeded in expanding the choice of research topics given in previous studies, as well as increasing the number of recommended supervisors and students. |

## 1. Introduction

Every student wishes to graduate on time because the length of time graduation can be used as an indicator of the student's quality. Based on research in [1-3], many students are still unable to complete their studies on time. One of the factors that influence the student's study period is the selection of supervisors [4]. Proper supervisor is a very important factor that will determine the success of the student [5]. Many studies have been done to create a recommendation system for selecting supervisors like in [6-12]. All of these researches only take student preferences without considering the supervisor's preferences, which we called a one-way recommendation system. According to Chen *et al.,* [13], when compared to a one-way recommendation system, a two-way

* Corresponding author.
E-mail address: annisa@apps.ipb.ac.id

recommendation system can satisfy both parties more because it takes into account the preferences of both parties [14].

A two-way recommendation system using skyline view queries [13], is used to consider the preferences of two parties and can answer queries needed by both parties. A two-way recommendation system is generally used for recommending job applicants (job seekers) to companies that are hiring or match making application. After determining the criteria that are taken into account by each party, then the job seeker input his/her job preference while the company input their preferences for the prospective employee. Skyline View Queries (SVQ) is a set of eight queries from both parties. The skyline view queries concept is based on the skyline query concept [15,16], which is a well-known method for obtaining a small set of dominant objects from large data sets known as skyline objects. An object is said dominant object if it is equally good in all dimensions and at least slightly better in one dimension [15].

Let $D = (D_1, \dots, D_d)$ be a multidimensional space and assume that the domain of each attribute $D_i (1 \le i \le d)$ is numeric. Let $S$ is a set of objects in $D$. Two objects $p, q \in S$. Notation $p \succ q$ means that $q$ is dominated by $p$ if for each attribute $D_i$, $p.D_i \ge q.D_i$, and there is at least one attribute $D_j (1 \le j \le d)$ where $p.D_j \ge q.D_j$. An object $p$ *is* said to be a skyline object if there does not exist $p' \succ p$ where $p'$ is another object in $S$. A set of skyline objects in $S$ denoted by $Sky(S)$.

In Ref. [17], a two-way recommendation system based on Skyline View Queries (SVQ) on Ref. [13] has been proposed for the selection of supervisors. Assume two parties in the supervisors' recommendation system are students and supervisors. Given a student data set $\mathcal{X}$ and supervisor data set $\mathcal{Y}$. A student has several criteria, we call it preference, for his supervisor and vice versa. Supervisors certainly want students who have high GPAs and course grades, while students want supervisors with the lowest average supervising time and have graduated many students. These preferences are modeled as predicates for students and supervisors. The preference of a student $X \in \mathcal{X}$ as a predicate $X.P$ in the space of the supervisor $\mathcal{Y}$, and vice versa. Student $X$ is interested in supervisor $Y \in \mathcal{Y}$ if $Y$ satisfies the predicate of $X$, $X.P(Y) = true$, and vice versa.

In contrast to the application of two-way recommendations in Ref. [13] where each party's interests are already available, in the selection of supervisors in Ref. [17], the interests of supervisors and students are not explicitly available. Most students and supervisors are rarely able to express their interests explicitly. To get the interests of supervisors and students who will be supervised (hereinafter referred to as students), this study uses historical data from students who have previously graduated (hereinafter referred to as alumni). Student interest in supervisors is obtained from the alumni's theses that were graduated on time. By using text mining, alumni theses are extracted to produce a collection of relevant topics. Furthermore, the students will select their preferred topics. Students who like topics that are similar to the alumni's thesis will be grouped into one cluster, and supervisors of the alumni can be assumed as supervisors that match the interest of students in the same cluster. Supervisor interest for students is obtained using the same method, using students' academic data and alumni's academic data who graduate on time. Students who have academic performance similar to alumni's will be grouped into one cluster, and it can be assumed that the supervisor of the alumni will be interested in these students.

Based on the above definitions, eight queries in SVQ for the supervisor's recommendation are defined as described in Table 1. The eight queries are applied symmetrically, for students and supervisors.

**Table 1**
The eight queries in Skyline View Queries for student ($X$) [13]

| Query number | Notation | Meaning |
|---|---|---|
| 1 | $V(X)$ | The view of $X$, $\{Y \in \mathcal{Y} \mid X.P(Y) = true\}$ |
| 2 | $IV(X)$ | The inverse view of $X$, $\{Y \in \mathcal{Y} \mid Y.P(A) = true\}$ |
| 3 | $SV(X)$ | The skyline view of $X$, $Sky(V(X))$ |
| 4 | $iSky(X)$ | The inverse skyline of $X$, $\{Y \mid X \in SV(Y)\}$ |
| 5 | $rSky(X)$ | The reciprocal skyline of $X$, $Sky(IV(X))$ |
| 6 | $MV(X)$ | The mutual view of $X$, $\{Y \mid Y \in V(X) \wedge X \in V(Y)\}$ |
| 7 | $SMV(X)$ | The skyline mutual view of $X$, $Sky(MV(X))$ |
| 8 | $SIS(X)$ | The skyline of inverse skyline of $X$, $Sky(iSky(X))$ |

The recommendation system proposed in Ref. [17] has major limitations. Some students and supervisors are still not recommended to anyone. This is due to the large number of students and supervisors who are not the skyline objects. Moreover, since Ref. [17] only used the title of the thesis to capture the student's interest, the research topics to be offered to students are small. To improve the recommendation model, in this research we use the k-Skyband Query concept [18] in SVQ instead of the skyline query. We also use the abstract of the thesis in addition to the thesis title.

The k-Skyband Query can be seen as an extension of traditional skylines, with a substantial difference: the parameter k allows to specify the maximum number of points that can dominate an element of the skyline. The usage of k allows Skyband Query to be more flexible concerning classical skylines. The k-skyband contains only the points that are dominated by at most $k$ points [18].

Using the k-skyband, students, and supervisors who are not in the skyline in Ref. [17] can be selected as skyband objects, and recommended to the target party. The contributions of this paper are summarized below:

 i. We have modified the concept of skyline view queries using the k-skyband query concept, and call it the k-skyband view queries.
 ii. We have improved the recommendation result and the number of research topics offered to students in [17].
 iii. We have conducted some experiments to prove the superiority of our proposed method.

The rest of this paper is organized as follows. In Section 2 we present the k-skyband view queries. In Section 3, we present the result of our research. Finally, in Section 4, we conclude the paper and show directions for our future works.

## 2. Methodology
### 2.1 Research Data and Environment

In this study, we use historical data of students from the Department of Computer Science, IPB University from 2010 to 2018. The length of study time, student academic data, and student thesis data were all included in the historical data. The computer used has an Intel Core i7-7700HQ Processor, 8GB of RAM, an NVIDIA GeForce GTX 1050 Graphics Processing Unit (GPU) with 4GB of memory (VRAM), and a 64-bit Windows 10 Home operating system. Jupyter Notebook, Visual Studio Code, Python 3.7, and an internet browser were all used in this research.

## 2.2 Research Stages

We reused the stages in Ref. [17], with modifications in the step to get student interest and recommendation calculation stages (marked with bold border pattern fill square) in Figure 1.
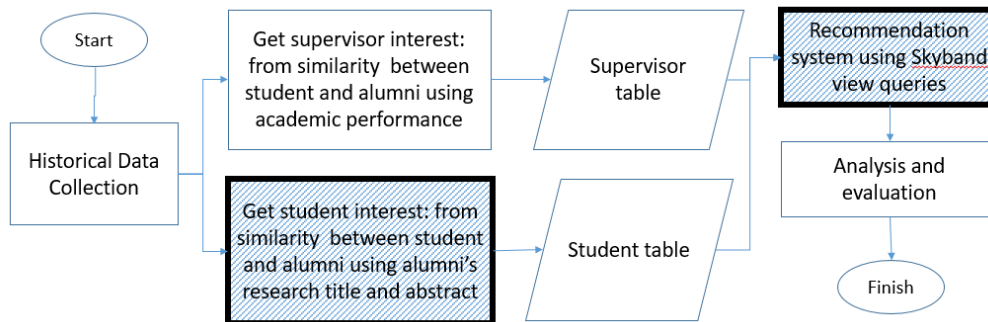


**Fig. 1.** Research Stages

For simplicity, we use Table 2 and Table 3 as examples of student and supervisor datasets.

**Table 2**
The student dataset

| Student | GPA | Course grade | Interest |
|---------|-----|--------------|----------|
| M1 | 3.5 | A | ? |
| M2 | 3.9 | B | ? |
| M3 | 3 | AB | ? |
| M4 | 3.3 | BC | ? |

**Table 3**
The supervisor dataset

| Supervisor | Average of student's study time (in years) | Number of supervised student | Interest |
|------------|---------------------------------------------|------------------------------|----------|
| D1 | 2 | 8 | ? |
| D2 | 2,5 | 12 | ? |
| D3 | 3 | 10 | ? |

## 2.2.1 Obtain student interest

To fill in the student interest in Table 2, titles and abstracts from alumni's theses are extracted to obtain relevant topics. Figure 2 depicts the topic extraction process. Each word in the title and abstract of the alumni thesis will be labeled based on the type of word in the word labeling section with Part of Speech (POS) tagging. For example, the word "I" is labeled PRP because it is a first-person pronoun, whereas "data" is labeled NN because it is a noun. The labeling guidelines in this research are based on the research provided and adhere to the guidelines used in Ref. [19]. After being labeled, the thesis title and abstract will be divided into sentences of two to three words using n-grams. The combined labels for each word contained in an n-gram will form a label pattern. For example, the sentence "my data" will have the NN PRP label pattern. The appearance of the label and sentence pattern is then counted.
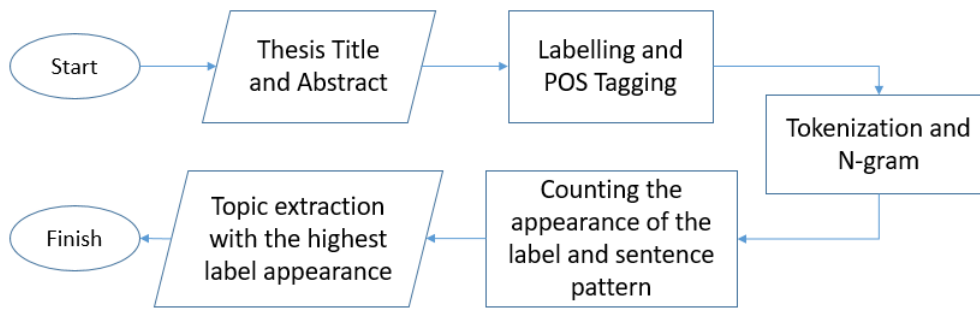
**Fig. 2.** Topic extraction Process

This collection of relevant topics is then given to students. Students choose the topics they are interested and we form the vector of student interest topics. We also form the same vector for alumni based on their research. Using the Jaccard similarity value [20], we calculate the similarity value between student vectors and alumni vectors. We assumed that the higher the similarity value is, the similar student and alumni interest is. It means the student is also assumed to have similar interests to the alumni's supervisor. Furthermore, the names of the supervisors obtained can be filled in the student interest column in Table 2.

### 2.2.2 Obtain supervisor interest

To fill in the supervisor interest in Table 3, we use academic data from alumni and students. Student academic data is collected to form a student academic data vector. The same process is also applied to alumni academic data. From these vectors, the clustering step is then carried out. Alumni and students who are in one cluster are assumed to have similar performance so supervisors from alumni will be interested in students in the same cluster. Furthermore, supervisor interest in Table 3 can be filled in with the names of students who are in the same cluster as the alumni supervised by the supervisor.

### 2.2.3 Skyband View Queries

In Ref. [17], the process of obtaining all skyline view query results is carried out in several stages as described in Figure 3. Modifications are made by applying the k-skyband concept to SV, iSKy, SIS, rSky, and SMV as shown in Figure 4.
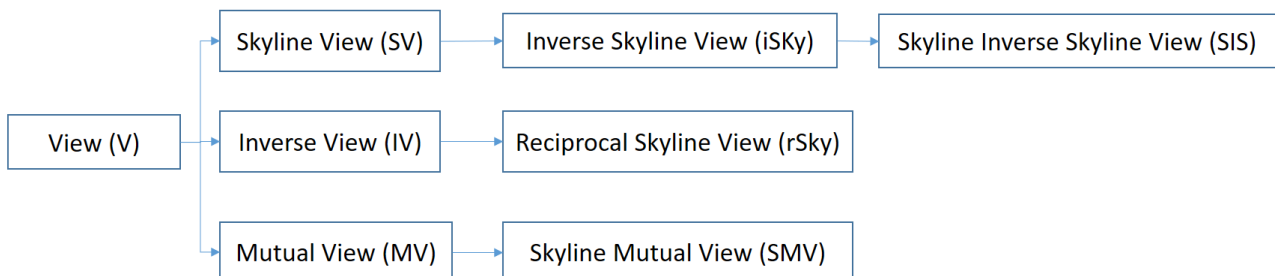


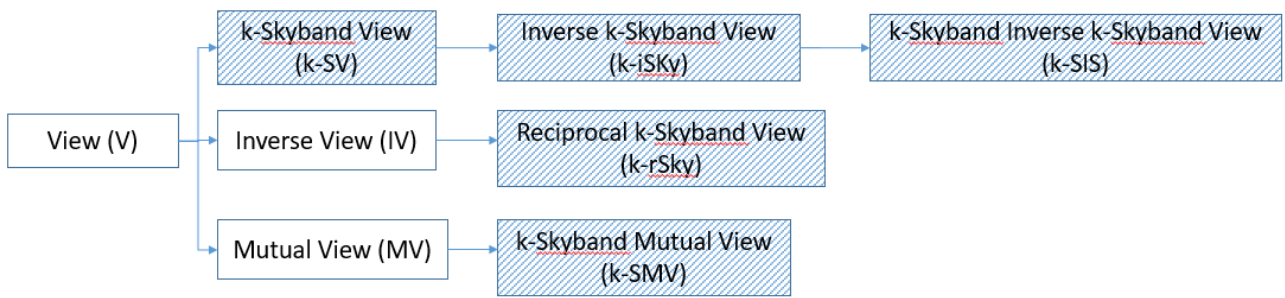**Fig. 3.** The Skyline view queries

**Fig. 4.** The k-Skyband view queries

The k-skyband concept makes it possible for all supervisors and students to be recommended. The application of the k-skyband concept is illustrated in Figure 5. Assume there are 10 students, M1-M10, plotted in two dimensions, student grade, and GPA. The higher the value of the two dimensions, the better. The skyline query will only return the most dominating students/supervisors (Figure 5(a)). By applying the k-skyband and adding $k$ values, we can add students who are not a skyline but still dominate other students, as shown in Figure 5(b) ($k$=1) and Figure 5(c) ($k$=2).
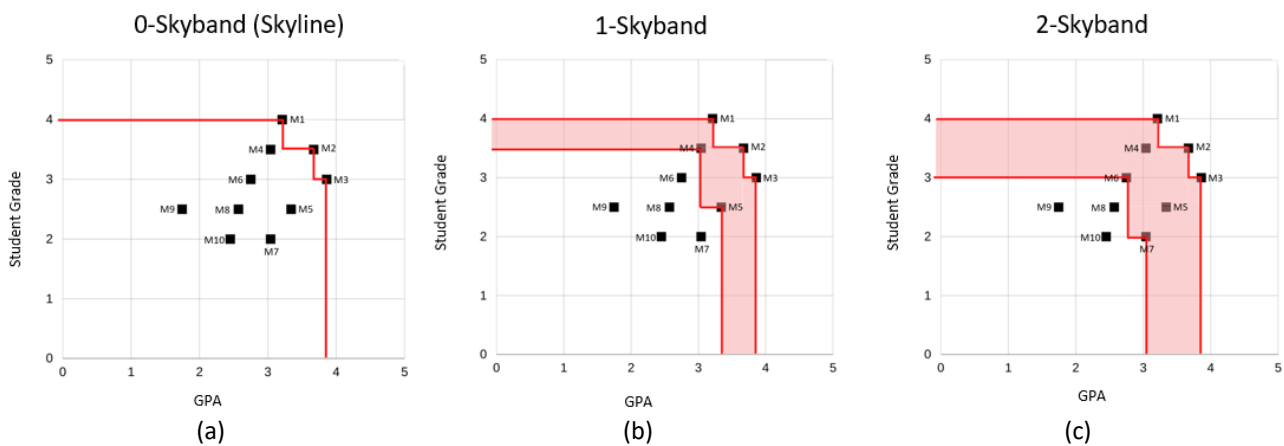


**Fig. 5.** (a) the Skyline query (b) the 1-Skyband query (c) the 2-Skyband query

*2.2.4 Analysis and evaluation*

The recommendations generated by the system will have four possibilities values [21] as shown in Table 4.

**Table 4**
The four possibilities value of recommendation result [21]

|  | Recommended | Not Recommended |
|---|---|---|
| Used | *True-Positive* (TP) | *False-Negative* (FN) |
| Not Used | *False-Positive* (FP) | *True-Negative* (TN) |

To evaluate the recommendation results, in this research, we use three metrics used in Ref. [21]. The first metric is precision, to assess the accuracy of the recommendation results. The second metric is recall, to assess the fulfillment of the user needs against the recommendation result. The last metric is the FP rate, to assess how many recommendations are not selected. The calculation of the precision, recall, and FR rate can be seen in Eq. (1), (2), and (3). The evaluation is also performed by comparing the results of the k-skyband view query recommendations with the recommendation result in Ref. [17].

$$Precision = \frac{TP}{TP + FP} \tag{1}$$

$$Recall = \frac{TP}{TP + FN} \tag{2}$$

$$False\ Positive\ Rate\ (FP\ Rate) = \frac{FP}{FP + TN} \tag{3}$$

## 3. Results
### 3.1 Topic Extraction Results

Topic extraction from the title data and abstracts obtained relevant thesis topics for students to choose from. Words in Indonesian are labels using POS Tagging with the tagger provided by Ref. [19]. For English words, we use the tagger provided by the Natural Language Tool Kit (NLTK) library. After tagging the thesis title and abstract data, it is divided into sentences using n-grams. The number of occurrences for each extracted sentence is calculated using n-grams after obtaining the number of occurrences of the label pattern. After that, topic extraction is performed. The results of the topic extraction stage were 258 topics that could be chosen by students, 5 times more than those produced in Ref. [17].

### 3.2 Analysis and Evaluation of the k-skyband View Queries

The recommendation results for SIS, rSky, and SMV in skyline view queries and k-SIS, k-rSky, and k-SMV in k-skyband view queries are shown in Figure 6(a) (from the student's perspectives) and 6(b) (from the supervisor's perspectives). Based on both student and supervisor perspectives, the number of recommendations resulting from k-skyband view queries, in general, is larger than skyline view queries, so it ensures that all lecturers and students are recommended.
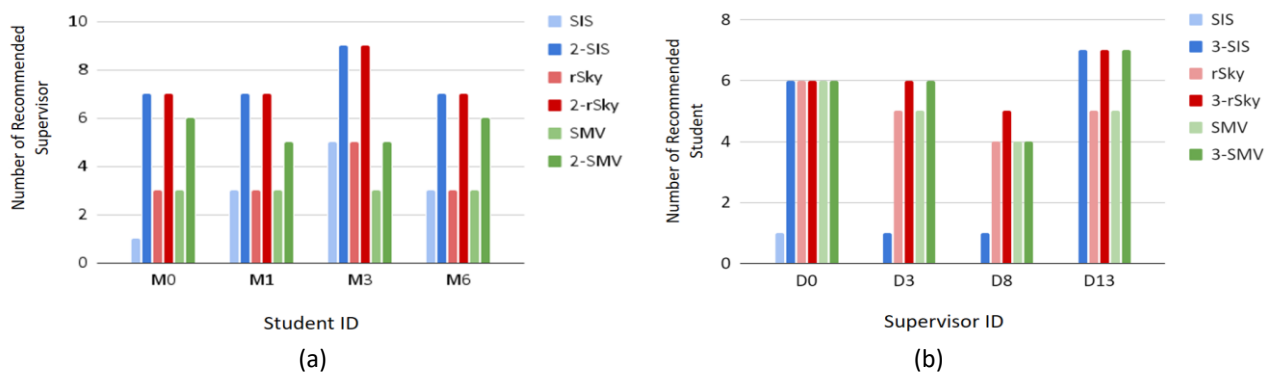


**Fig. 6.** (a) the number of recommendation results for skyline view and skyband view queries based on the student's perspective (b) supervisor's perspective

Figure 7(a) and (b) depict the effect of $k$ on the number of supervisors and students recommended for query SIS, rSky, and SMV. As can be seen, increasing the value of $k$ will increase the recommended number of students and supervisors. Nonetheless, it appears that the higher $k$, the higher the number of recommendation results is.
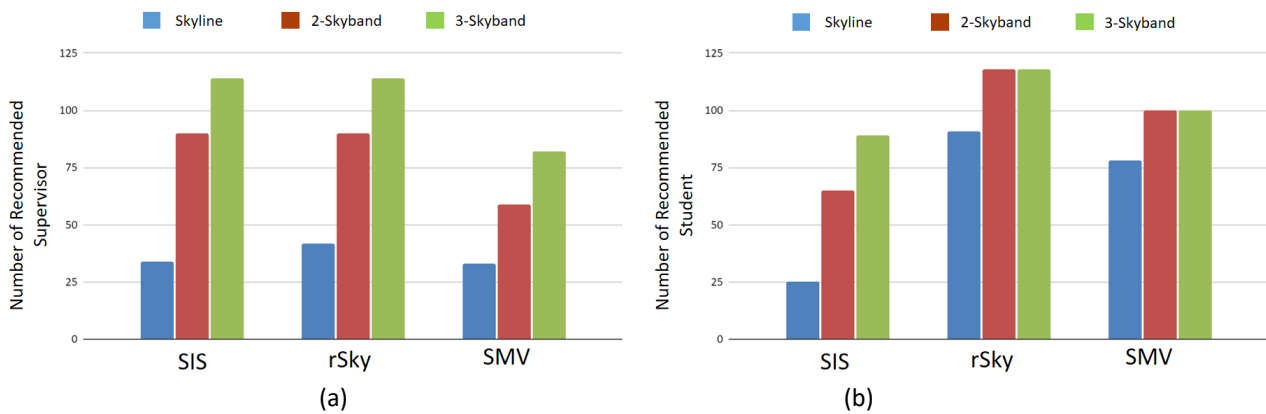
**Fig. 7.** (a) the effect of the value of $k$ on the number of recommended supervisors and (b) on the number of recommended students

The processing time is also affected by the higher $k$ value. Figure 8 shows the effect of increasing $k$ values on processing time. The increase in the processing time is caused by k more skyline iterations should be done when k is increased. However, a significant increase in processing time occurs at low k values, at higher k values the processing time does not show a significant increase.
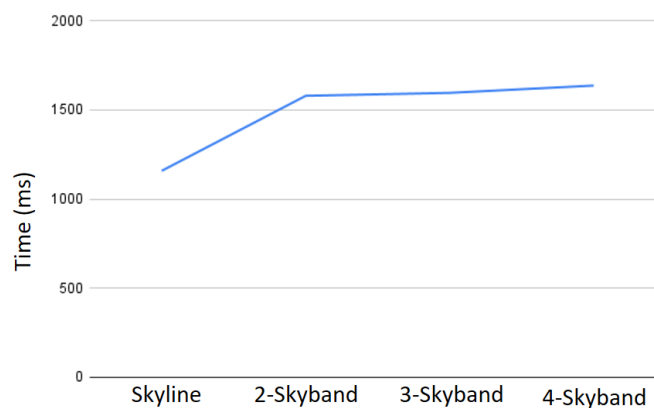


**Fig. 8.** The effect of the value of $k$ on processing time

The higher the value of $k$ used in skyband view queries, the more students and supervisors who are less dominant will be recommended. The results of recommendations can be prioritized based on the number of $k$ used. Students who are recommended to supervisors with 1-skyband results are given high priority. If the supervisor cannot accept him/her, the student will be recommended to the 2-skyband supervisors' results and so on.

Next, precision, recall, and FP rate are calculated to evaluate the results of the supervisor's recommendations. Figure 9 shows the results of the three metrics' calculations.

As shown in Figure 9(a), adding more k to k-skyband view queries does not always increase the precision of the recommendation results. This is because the results of the recommendations given are not necessarily chosen by the student. However, increasing the value of k can increase the recall value as shown in Figure 9(b). This is because the more recommendations given, the more supervisors are by the wishes of students so more supervisors will be selected. However, the more recommendations given, the higher the FP Rate will be, where more recommendation results will not be selected.
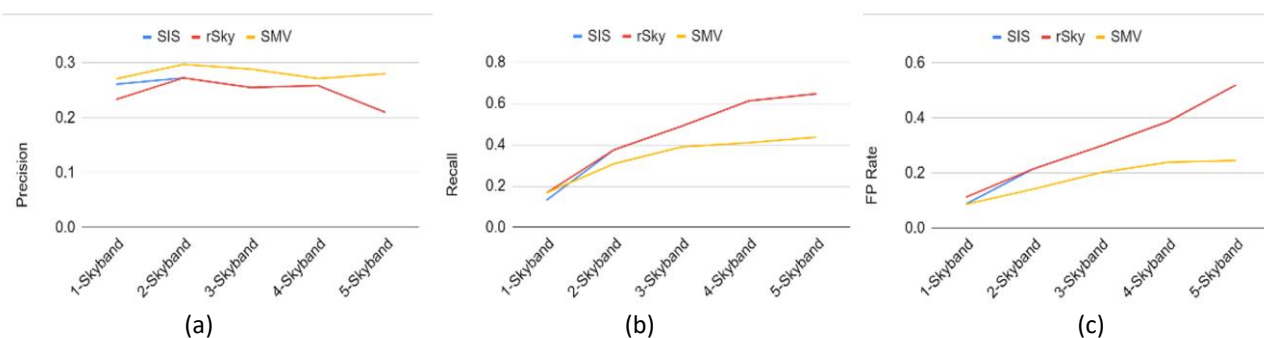
312

**Fig. 9.** (a) Precision from the recommendation of the supervisor (b) Recall from the recommendation of the supervisor (c) FP Rate from the recommendation of the supervisor

## 4. Conclusions

The most dominant supervisors and students are recommended by the skyline view query recommendation model. However, there are still many students and supervisors who are not recommended because they are not the most dominant. Applying the k-skyband concept to existing models can increase the number of recommended students and supervisors. This study was successful in providing supervisor recommendations to all students as well as alternative options for supervisors based on the value of k on the k-skyband. This study was also successful in broadening the topics available to students. Topics are expanded by including abstract data and extracted the abstract data using n-grams. The most relevant topics are chosen by sorting the number of occurrences of sentences and the number of occurrences of label patterns given with POS tagging of all thesis and final project documents. This study was successful in increasing the number of topics available for students. Adding a k value to k-skyband view queries can improve user satisfaction by allowing students to select supervisors of interest. However, as the value of k increases, there will be more supervisors who are recommended but not chosen by students, allowing the FP rate to rise as a result of the recommendations.

**References**
[1] Gnjidic, Danijela, Narelle da Costa, and Nial J. Wheate. "Potential factors that can affect the performance of undergraduate pharmacy research students: a descriptive study." *BMC Medical Education* 23, no. 1 (2023): 1-9. https://doi.org/10.1186/s12909-023-04018-5
[2] Srinadi, I. G. A. M., and Desak Putu Eka Nilakusmawati. "Analisis Waktu Kelulusan Mahasiswa Fmipa Universitas Udayana Dan Faktor-Faktor Yang Memengaruhinya." *E-Jurnal Matematika Udayana University* 9, no. 3 (2020): 205-212. https://doi.org/10.24843/MTK.2020.v09.i03.p300
[3] Nofratama, Fadli, Hasrul Hasrul, Henni Muchtar, and Susi Fitria Dewi. "Kendala Keterlambatan Penyelesaian Studi Mahasiswa PPKn FIS Universitas Negeri Padang." *Journal of Education, Cultural and Politics* 2, no. 2 (2022): 185-191. https://doi.org/10.24036/jecco.v2i2.106
[4] Yuniar, Dina, Heti Mulyati, and Eko Ruddy Cahyadi. "Faktor-Faktor Yang Mempengaruhi Penyelesaian Masa Studi Program Pascasarjana Di Institut Pertanian Bogor." *Jurnal Akuntabilitas Manajemen Pendidikan* 7, no. 2 (2019): 133-147. https://doi.org/10.21831/amp.v7i2.25084
[5] Malik, Mohd Azry Abdul, Nur Izzatulsyimah Madzuki, Nur Syahidah Shahnirul Hizam, Nuramanina Husna Shamsul Kamal, Nur Syaliza Hanim Che Yusof, Mohd Faiez Suhaimin, and Siti Nurani Zulkifli. "Teachers' Readiness and Practices in School-Based Assessment Implementation: Primary Education in Malaysia." *International Journal of Advanced Research in Future Ready Learning and Education* 23, no. 1 (2021): 1-9.

[6]     Rismanto, Ridwan, Arie Rachmad Syulistyo, and Bebby Pramudya Citra Agusta. "Research supervisor recommendation system based on topic conformity." *International Journal of Modern Education and Computer Science* 12, no. 1 (2020): 26. https://doi.org/10.5815/ijmecs.2020.01.04

[7]     Arisetiawan, Anak Agung Bagus, Indriati Indriati, and Dian Eka Ratnawati. "Sistem Rekomendasi Dosen Pembimbing Berdasarkan Dokumen Judul Skripsi di Bidang Komputasi Cerdas Menggunakan Metode BM25." *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer* 3, no. 6 (2019): 5832-5836.

[8]     Yulianti, Wilda, Budi Sutomo, and Andreas Perdana. "Sistem Rekomendasi Pemilihan Dosen Pembimbing Skripsi Menggunakan Metode Weighted Product (Studi Kasus STMIK Dharma Wacana Metro)." *I-Robot* 2 (2018): 341338. https://doi.org/10.53514/ir.v2i1.79

[9]     Asfi, Marsani, and Nopi Fitrianingsih. "Implementasi Algoritma Naive Bayes Classifier sebagai Sistem Rekomendasi Pembimbing Skripsi." *InfoTekJar: Jurnal Nasional Informatika dan Teknologi Jaringan* 5, no. 1 (2020): 44-50.

[10]    Abdullah, Asrul, and Menur Wahyu Pangestika. "Perancangan Sistem Pendukung Keputusan Dalam Pemilihan Dosen Pembimbing Skripsi Berdasarkan Minat Mahasiswa dengan Metode AHP (Analytical Hierarchy Process) di Universitas Muhammadiyah Pontianak." *JEPIN (Jurnal Edukasi dan Penelitian Informatika)* 4, no. 2 (2018): 184-191. https://doi.org/10.26418/jp.v4i2.27651

[11]    Kazakovtsev, Vladimir, Svyatoslav Oreshin, Alexey Serdyukov, Egor Krasheninnikov, Sergey Muravyov, Albert Bezvinnyi, Alexander Panfilov et al. "Recommender system for an academic supervisor with a matrix normalization approach." In *Proceedings of the 2020 1st International Conference on Control, Robotics and Intelligent System*, pp. 84-87. 2020. https://doi.org/10.1145/3437802.3437817

[12]    Hasan, Mir Anamul, and Daniel G. Schwartz. "Recadvisor: Criteria-based ph. d. supervisor recommendation." In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pp. 1325-1328. 2018. https://doi.org/10.1145/3209978.3210178

[13]    Chen, Jian, Jin Huang, Bin Jiang, Jian Pei, and Jian Yin. "Recommendations for two-way selections using skyline view queries." *Knowledge and information systems* 34 (2013): 397-424. https://doi.org/10.1007/s10115-012-0489-6

[14]    Lin, Yi-Wen, En Tzu Wang, Chieh-Feng Chiang, and Arbee LP Chen. "Finding targets with the nearest favor neighbor and farthest disfavor neighbor by a skyline query." In *Proceedings of the 29th Annual ACM Symposium on Applied Computing*, pp. 821-826. 2014. https://doi.org/10.1145/2554850.2554863

[15]    Borzsony, Stephan, Donald Kossmann, and Konrad Stocker. "The skyline operator." In *Proceedings 17th international conference on data engineering*, pp. 421-430. IEEE, 2001.

[16]    Tiakas, Eleftherios, Apostolos N. Papadopoulos, and Yannis Manolopoulos. "Skyline queries: An introduction." In *2015 6th International Conference on Information, Intelligence, Systems and Applications (IISA)*, pp. 1-6. IEEE, 2015. https://doi.org/10.1109/IISA.2015.7388053

[17]    Sampurno, Global Ilham, Annisa Annisa, and Sony Hartono Wijaya. "Sistem Rekomendasi Dua Arah untuk Pemilihan Dosen Pembimbing Menggunakan Data Histori dan Skyline View Queries." *Jurnal Teknologi Informasi dan Ilmu Komputer* 9, no. 5 (2022): 1055-1064. https://doi.org/10.25126/jtiik.2022955458

[18]    Siddique, Md Anisuzzaman, Asif Zaman, and Yasuhiko Morimoto. "Selection of Important Sets by using K-Skyband Query for Sets." *International Journal of Advanced Computer Science and Applications* 9, no. 4 (2018). https://doi.org/10.14569/IJACSA.2018.090444

[19]    Dinakaramani, Arawinda, Fam Rashel, Andry Luthfi, and Ruli Manurung. "Designing an Indonesian part of speech tagset and manually tagged Indonesian corpus." In *2014 International Conference on Asian Language Processing (IALP)*, pp. 66-69. IEEE, 2014. https://doi.org/10.1109/IALP.2014.6973519

[20]    Chung, Neo Christopher, Błażej Miasojedow, Michał Startek, and Anna Gambin. "Jaccard/Tanimoto similarity test and estimation methods for biological presence-absence data." *BMC bioinformatics* 20, no. 15 (2019): 1-11. https://doi.org/10.1186/s12859-019-3118-5

[21]    Shani, Guy, and Asela Gunawardana. "Evaluating recommendation systems." *Recommender systems handbook* (2011): 257-297. https://doi.org/10.1007/978-0-387-85820-3_8