

Analisis Cluster Pada Pengelompokan Siswa Diktuk Bintara Polri TA. 2018/2019, SPN Singaraja - Polda Bali Menggunakan K-Means dan K-Harmonic Means

I Gede Aris Gunadi¹⁾, D. P. Singgih Putri²⁾, I. N. Mistanada³⁾, G. S. Prayoga⁴⁾, N. M. Y. D. Rahayu⁵⁾

¹⁾ Jurusan Pendidikan Fisika FMIPA / Prodi S2 Ilmu Komputer Universitas Pendidikan Ganesha

^{2,3,4,5)} Prodi S2 Ilmu Komputer Universitas Pendidikan Ganesha

¹⁾igedearisgunadi@undiksha.ac.id, ²⁾desysinggihputri@unud.ac.id, ³⁾nmistanada@gmail.com,

⁴⁾gama.pravoga666@gmail.com, ⁵⁾yeni.brt@gmail.com

ABSTRACT

The process of student grouping at Dikduk program in SPN (Sekolah polisi Negara) Singaraja is done with several processes. In the processes, there were 20 subject topics to be tested. The grouping of students is divided into five namely the technical functions, technical functions of the sabhara, the technical functions of the intelligence, the technical functions of the community organization and the technical function of the detective. There are difficulties in the grouping process which are caused by the large amount of data and the elements of the assessment carried out. In this study the grouping of students was done by using two clustering methods namely K-Means and K-Harmonic Means. The purpose of this study is to help the school in grouping large numbers of students into predetermined categories. This research was conducted involving 138 students. The students were divided into five groups based on 20 subjects as variables. The results showed that the K-Harmonic Means clustering method had a smaller number of data in the incorrect group compared to the K-Means method, and the index silhouette value on the K-Harmonic Means method was higher than the K-Means method.

Keywords: metode clustering, K-Means, K-Harmonic Means, Silhouette Indeks.

I. PENDAHULUAN

Pendidikan Pembentukan (Diktuk) Bintara Polri pada Sekolah Pendidikan Kepolisian Negara (SPN) Singaraja memiliki tujuan untuk membentuk insan Bhayangkara yang memiliki sikap, perilaku, pengetahuan, keterampilan, dengan kondisi fisik yang Samapta (siap siaga, siap sedia dan waspada) untuk melaksanakan tugas sebagai pemelihara keamanan dan ketertiban masyarakat, penegak hukum, pelindung, pengayom, dan pelayan masyarakat yang profesional, bermoral, modern dan unggul.

Pada Tahun Ajar (TA) 2018/2019 SPN Singaraja memiliki 138 siswa dengan nama angkatan Praja Raksaka Gautama (PRG 2). Siswa tersebut kemudian dibagi ke dalam lima kelompok yaitu fungsi teknis lantas, fungsi teknis sabhara, fungsi teknis intelegen, fungsi teknis binmas dan fungsi teknis reserse. Selama ini pengelompokan siswa dilakukan dengan menghitung berdasarkan kriteria maupun unsur-unsur yang terdapat pada fungsi teknis dan berdasarkan 20 mata pelajaran yang diampu. Dari 20 mata pelajaran tersebut terdiri dari pengantar yang berisi pengenalan lingkungan dan tradisi lembaga pendidikan dan pengarahan program. Mata pelajaran selanjutnya adalah kelompok mata pelajaran yang terdiri dari kepribadian, pengetahuan sosial, hukum, profesi teknis kepolisian dan jasmani.

Metode pembelajaran yang dilakukan pun memiliki beberapa macam yaitu metode tanya jawab, diskusi, ceramah, pemecahan masalah, penugasan, demonstrasi, latihan/ *drill*, simulasi dan *role play*. Jika metode-metode tersebut sudah dapat dilakukan maka para

siswa akan diberi nilai, penilaian tersebut didapat dari penilaian hasil belajar yang terdiri dari aspek akademik, mental kepribadian dan kondisi Samapta jasmani. Selanjutnya penilaian secara teknis diatur khusus dalam pedoman penilaian.

Dari pengelompokkan siswa yang dilakukan, dibutuhkan suatu metode yang mampu membantu dalam pengelompokkan dengan cepat dan sesuai. Sehingga metode yang digunakan dalam membagi siswa menjadi beberapa kelompok dengan menggunakan metode clustering yaitu metode *K-Means* dan *K-Harmonic Means*. Dari kedua metode tersebut diharapkan mampu mengelompokkan siswa menjadi lima kelompok fungsi teknis.

II. TINJAUAN PUSTAKA

2.1. K-Means

Metode K-Means Clustering (KM) merupakan metode klasterisasi secara partisi (*partitional clustering*). Hampir semua metode klasterisasi secara partisi didasarkan pada tujuan untuk mengoptimalkan nilai fungsi $f(x)$ sebagai *clustering criterion* dimana hal ini dapat dikatakan sebagai penerjemahan gagasan intuisi manusia terhadap suatu klaster kedalam suatu rumus matematis . KM tidak menjamin hasil klasterisasi yang unik karena metode ini dapat menghasilkan hasil klaster yang berbeda tergantung dari posisi inisialisasi klaster awal (S. S. Khan & Ahmad, 2004)

Berikut ini gambaran dari K-Means. Misal $X = \{x_i | i = 1, \dots, n\}$ merupakan suatu himpunan n titik berdimensi d yang akan diklasterkan kedalam K klaster $C = \{c_k | k = 1, \dots, K\}$. Metode K-Means menemukan suatu partisi/klaster sedemikian hingga nilai *squared error* antara titik tengah (*mean*) dari suatu klaster ke semua titik data klaster tersebut merupakan nilai minimum (Kumar Dehariya, 2010). Misalkan μ_k adalah rata-rata dari klaster c_k yang didapat dari persamaan 1.

$$\mu_k = \frac{1}{n_k} \sum_{x_i \in c_k} x_i \dots\dots\dots (1)$$

Dimana n_k merupakan jumlah elemen pada c_k . Squared error antara μ_k dan seluruh data pada klaster c_k didasarkan pada jarak Euclidean antara titik yang ada dengan pusat klasternya, *squared error* tersebut didefinisikan seperti pada persamaan 2.

$$J(c_k) = \sum_{x_i \in c_k} \|x_i - \mu_k\|^2 \dots\dots\dots (2)$$

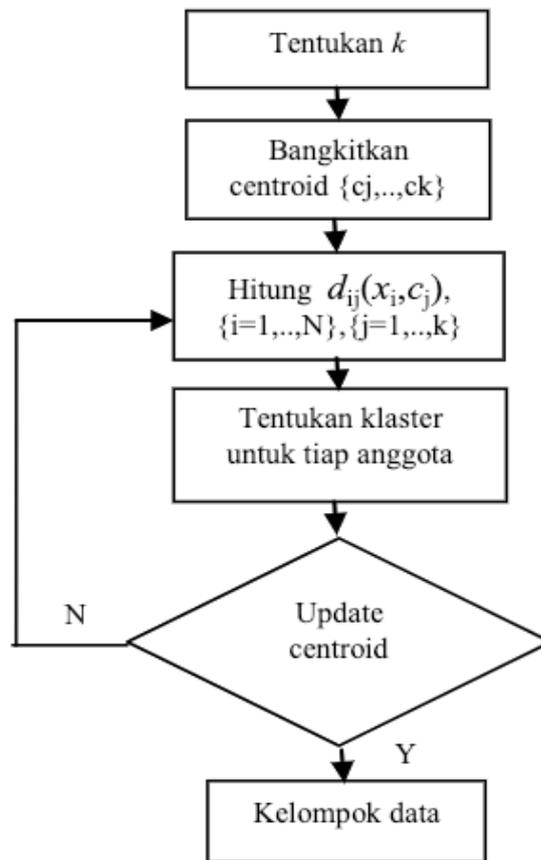
Fungsi tujuan (*objective function*) dari klasterisasi dengan K-Means adalah meminimkan *total squared error* dari seluruh klaster. Fungsi tujuan ini juga disebut sebagai *clustering criterion* dan juga sebagai *cost function* (S. S. Khan & Ahmad, 2004) dalam penemuan solusi optimal. Adapun formula dari tujuan ini seperti persamaan 3.

$$J(C) = \sum_{k=1}^K \sum_{x_i \in c_k} \|x_i - \mu_k\|^2 \dots\dots\dots (3)$$

Solusi pada metode K-Means adalah terbentuknya klaster-klaster dengan nilai $J(C)$ yang minimum. Berikut adalah algoritma metode K-Means (Cebeci & Yildiz, 2015) Inisialisasi K titik pusat klaster awal secara acak

1. Klasterkan setiap obyek yang ada sesuai jarak terdekat ke pusat klaster yang ada
2. Perbaiki nilai semua pusat klaster
3. Ulangi langkah 2 dan 3 sampai nilai semua pusat klaster tidak ada perubahan.

Diagram alir algoritma K-Means ditunjukkan pada Gambar 1.



Gambar 1. Diagram Alir Prosedur K-Means

Dalam beberapa penelitian yang telah dilakukan, K-Means memiliki kekuatan komputasi yang sangat baik, (Cebeci & Yildiz, 2015) melakukan penelitian yang mencoba untuk mengkomparasi kecepatan antara K-Means clustering dengan Fuzzy C Mean. Dalam penelitian tersebut dinyatakan bahwa K-Means clustering memiliki kecepatan yang jauh lebih baik dibandingkan Fuzzy C Means Clustering. Namun dari segi akurasi dinyatakan Fuzzy C means clustering masih lebih baik. Dalam penelitian lain, (Z. Khan, Ni, Fan, & Shi, 2017), juga dilakukan uji coba terkait dengan kekuatan algoritma K-mean, yang dibandingkan dengan K-medoids dan fuzzy c means, untuk proses clustering pada data citra. Dalam penelitian ini dinyatakan salah satu permasalahan algoritma K mean adalah bagaimana menentukan parameter awal, yang dijadikan acuan sebagai titik pusat clustering. Dalam penelitian ditemukan bahwa algoritma K mean memiliki *performance* yang lebih baik dari k medoids maupun fuzzy c means. Untuk permasalahan klustering data citra. Pada kasus tertentu K Means juga dianggap memiliki kelemahan

2.2. K-Harmonic Means (KHM)

Algoritma K-Harmonic Means merupakan pengembangan dari K-Means yang memperbaiki kekurangan dari K-Means dengan menggunakan fungsi obyektif yang didapatkan dengan cara meminimalisasi rata-rata harmonik dari jarak seluruh data dengan tiap centroid. Dari hasil penelitian menunjukkan bahwa K-Harmonic Means tidak sensitif terhadap inisialisasi centroid dan secara signifikan meningkatkan kualitas klusterisasi dibandingkan dengan K-Means (Zhang, Hsu, & Dayal, 1999). Langkah-langkah K-

Harmonic Means sebagai berikut (Wicaksana & Widiartha, 2012) : pertama inialisasi posisi titik pusat kluster awal secara random. Kemudian jika p adalah input parameter dan biasanya nilai $p \geq 2$. hitung nilai fungsi persamaan 4.

$$KHM(X, C) = \sum_{i=1}^N \frac{K}{\sum_{l=1}^K \frac{1}{\|x_i - c_l\|^p}} \dots\dots\dots (4)$$

Selanjutnya untuk setiap data x_i , hitung nilai keanggotaan $m(c_l/x_i)$ untuk setiap titik pusat kluster c_l berdasarkan persamaan 5.

$$m(c_l | x_i) = \frac{\|x_i - c_l\|^{-p-2}}{\sum_{l=1}^k \|x_i - c_l\|^{-p-2}} \dots\dots\dots (5)$$

Untuk setiap data x_i , hitung nilai bobot $w(x_i)$ berdasarkan persamaan 6.

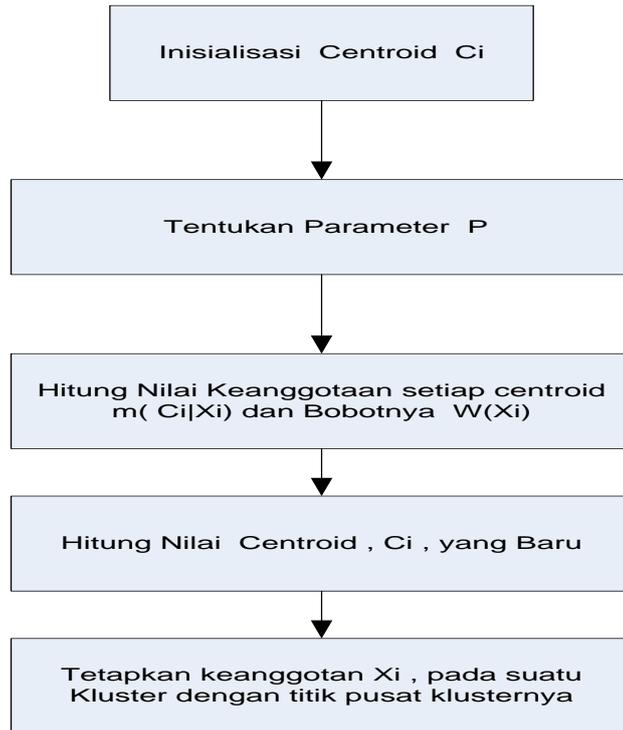
$$w(x_i) = \frac{\sum_{l=1}^K \|x_i - c_l\|^{-p-2}}{\left(\sum_{l=1}^K \|x_i - c_l\|^{-p}\right)^2} \dots\dots\dots (6)$$

Untuk setiap titik pusat x_i , ulang kembali perhitungan untuk posisi titik pusat kluster dari semua data berdasarkan nilai keanggotaan dan bobot yang dimiliki tiap data. Penentuan posisi titik pusat ini berdasarkan persamaan 7.

$$c_l = \frac{\sum_{i=1}^N m(c_l | x_i).w(x_i).x_i}{\sum_{i=1}^N m(c_l | x_i).w(x_i)} \dots\dots\dots (7)$$

Ulangi perhitungan nilai fungsi tujuan sampai *update* centroid sampai mendapatkan nilai fungsi tujuan yang tidak terdapat perubahan atau kurang dari ambang batas (Sani, 2018). Tetapkan keanggotaan data x_i pada suatu kluster dengan titik pusat kluster c_j sesuai dengan nilai keanggotaan x_i terhadap c_j .

Tahapan K-Harmonic Means disajikan pada Gambar 2.



Gambar 2. Diagram Prosedur K-Harmonic Means

2.3. Silhouette Indeks

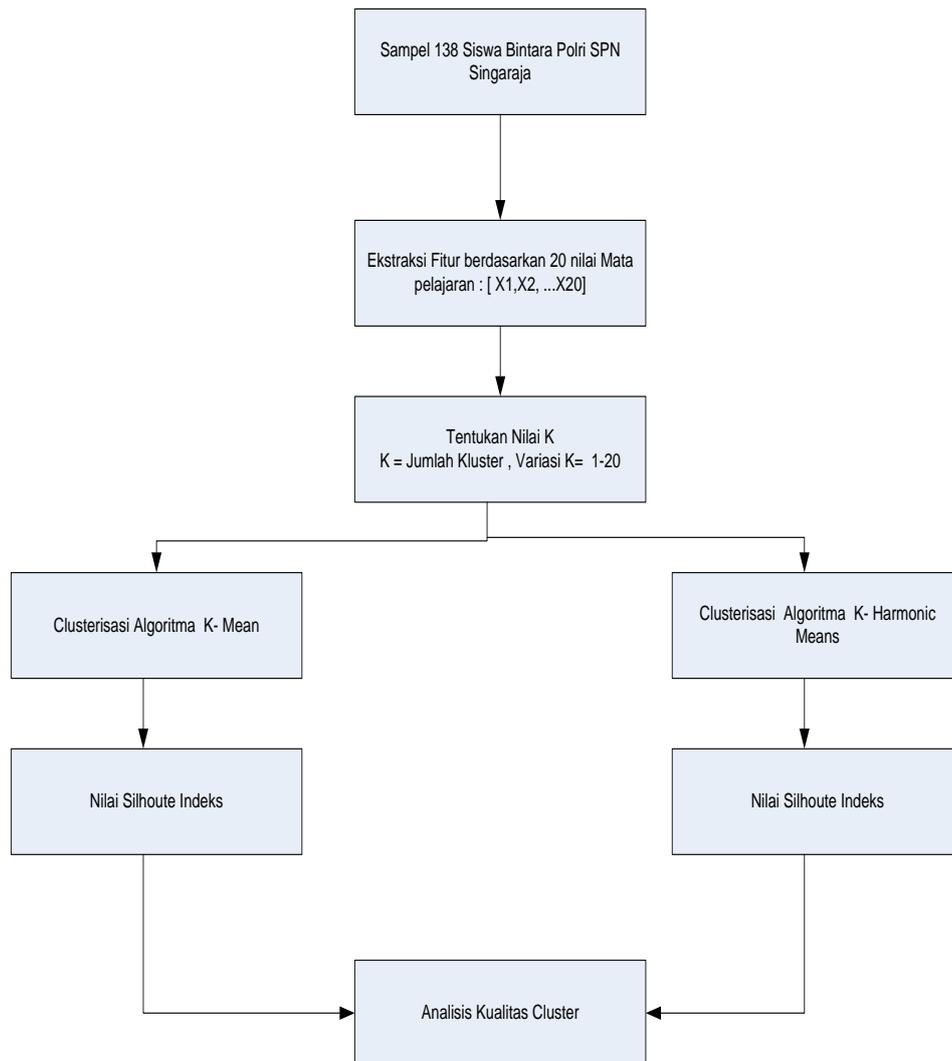
Silhouette Index (SI) digunakan untuk memvalidasi sebuah data, kluster tunggal, atau bahkan keseluruhan klaster. Metode ini banyak digunakan untuk memvalidasi kluster yang menggabungkan nilai kohesi dan separasi. Rentang nilai SI adalah -1 hingga +1. Nilai SI mendekati 1 menunjukkan bahwa data tersebut tidak tepat berada pada kluster tersebut. SI bernilai 0 atau mendekati 0 maka posisi data berada pada perbatasan dua klaster (Handoyo et al., 2014), (Gita Premashanti Trayasiwi, 2017), (Izzah & Hayatin, 2013), (Mukarromah, Martha, & Ilhamsyah, 2015).

Berikut ini adalah beberapa rumus yang digunakan untuk menghitung *Silhouette Indeks*, persamaan 8.

$$\begin{aligned}
 a_i^j &= \frac{1}{m_j - 1} \sum_{r=1}^{m_j} d(x_i^j, x_r^j) \\
 i &= 1, 2, \dots, m_j \\
 b_i^j &= \min \left\{ \frac{1}{m_n} \sum_{r=1}^{m_n} d(x_i^j, x_r^n) \right\}, i \\
 &= 1, 2, \dots, m_n \\
 SI_i^j &= \frac{b_i^j - a_i^j}{\max\{a_i^j, b_i^j\}} \\
 SI_j &= \frac{1}{m_j} \sum_{i=1}^{m_j} SI_i^j \\
 SI &= \frac{1}{k} \sum_{j=1}^k SI_j
 \end{aligned} \tag{8}$$

III. METODELOGI PENELITIAN

Penelitian ini bertujuan untuk melihat kualitas hasil *clusteing* yang dihasilkan oleh K-Means dan K -Harmonik Mean (KHM), dengan menggunakan paramater *silhouette indeks*, seperti yang dijelaskan pada Persamaan 8. Pengujian menggunakan 138 data calon bintanga Polri, dengan menggunakan 20 fitur nilai mata pelajaran. Pengujian dilakukan dengan menggunakan variasi nilai K dari 2 sampai 11. Pada setiap nilai K, dibandingkan nilai *silhouette indeks* antara algoritma K-Mean dan K-Harmonik Means. Secara umum ilustrasi tahapan penelitian dinyatakan pada Gambar 3.



Gambar 3 . Tahapan Penelitian

IV. HASIL DAN PEMBAHASAN

Dalam penelitian ini dipergunakan 138 data nilai siswa Diktuk Bintang Polri Tahun Ajar 2018/2019 yang masing-masing terdiri dari 20 mata pelajaran yang nantinya dijadikan variabel data x_1, x_2, \dots, x_{20} .

3. 1 Metode K-Means

Pengelompokkan menggunakan K-Means menghasilkan hasil seperti pada Tabel 1.

Tabel 1. Uji Validitas Cluster

Nilai K	Uji Validitas Cluster		
	Tepat	Berada di perbatasan	Tidak tepat
2	116	0	22
3	124	0	14
4	125	0	13
5	116	0	22
6	120	0	18
7	120	0	18
8	124	0	14
9	115	0	23
10	127	0	11
11	117	0	21

Pada Tabel 1 diujikan nilai K yang bervariasi, dengan variasi dari 2 sampai 11. Sebagai contoh untuk K sama dengan 2, dari 138 siswa dikelompokkan menjadi 2 kluster dan ditemukan 116 siswa yang memang berada pada kelompok yang bersesuaian. Adapun kategori kesesuaian adalah dilihat tingkat perbedaan nilai, dari 20 mata pelajaran. Dapat dinyatakan pada kasus K sama dengan 2 terdapat 116 siswa yang memiliki kemiripan nilai, dan telah dipetakan dalam kluster yang benar. Sebaliknya sisanya sebanyak 22 siswa berada pada kluster yang terdapat perbedaan nilai, dapat dikatakan 22 siswa berada pada cluster yang tidak tepat.

Berdasarkan persamaan 8, dilakukan perhitungan nilai SI (*Silhouette Indeks*), pada setiap variasi nilai K. Hasil perhitungan disajikan pada Tabel 2.

Tabel 2 Nilai SI dari Masing – Masing Kelompok

Nilai K	Nilai SI											
	rata2	Kel-1	Kel-2	Kel-3	Kel-4	Kel-5	Kel-6	Kel-7	Kel-8	Kel-9	Kel-10	Kel-11
2	0.212	0.358	0.066									
3	0.166	0.207	0.261	0.030								
4	0.144	0.208	0.028	0.261	0.080							
5	0.136	0.015	0.167	0.219	0.267	0.015						
6	0.136	0.103	0.275	0.043	0.205	0.194	-0.003					
7	0.102	0.059	0.111	0.002	0.112	0.189	0.051	0.190				
8	0.143	0.092	0.098	0.260	0.185	0.118	0.007	0.213	0.171			
9	0.114	0.121	-0.002	0.195	0.183	0.089	0.106	0.210	-0.011	0.131		
10	0.121	0.147	0.026	0.211	0.175	0.086	0.124	0.076	0.114	0.171	0.079	
11	0.122	0.073	0.033	0.243	0.213	0.077	0.172	0.099	0.185	0.104	0.079	0.059

Jika dilihat dari jumlah data yang berada pada kelompok yang tidak tepat, maka nilai K=10, K=4 dan K=3 menduduki peringkat 3 teratas. Namun apabila dilihat dari rata-rata nilai SI maka K=2, K=3 dan K=4 yang menduduki peringkat 3 teratas. Dengan demikian apabila menghendaki data yang berada pada kelompok yang tidak tepat berjumlah sedikit dengan nilai rata-rata SI yang tinggi, maka nilai K=3 dan K=4 layak untuk dijadikan pilihan.

3.2 Metode KHM

Menggunakan metode KHM menghasilkan hasil seperti pada Tabel 3.

Tabel 3. Uji Validitas Cluster dengan Metode KHM

Nilai K	Uji Validitas Cluster		
	Tepat	Berada di perbatasan	Salah
2	129	0	9
3	127	0	11
4	124	0	14
5	123	0	15
6	113	0	25
7	114	0	24
8	115	0	23
9	112	0	26
10	107	0	31
11	93	0	45

Dengan nilai SI untuk masing-masing nilai K seperti pada Tabel 4.

Tabel 4. Nilai SI dari Masing – masing Kelompok

Nilai K	Nilai SI											
	rata2	x1	x2	x3	x4	x5	x6	x7	x8	x9	x10	x11
2	0.229	0.098	0.359									
3	0.139	0.102	0.263	0.051								
4	0.136	0.264	0.020	0.187	0.074							
5	0.159	0.270	0.036	0.043	0.226	0.219						
6	0.131	0.083	0.022	-0.006	0.209	0.221	0.260					
7	0.112	0.002	0.051	0.048	0.207	0.185	0.230	0.064				
8	0.100	-0.025	0.060	0.212	0.110	0.043	0.016	0.210	0.176			
9	0.079	-0.005	0.024	0.170	0.066	0.203	0.010	0.205	0.089	-0.048		
10	0.072	0.078	0.036	0.065	0.186	0.146	0.066	0.013	-0.024	0.131	0.018	
11	0.120	0.154	-0.060	-0.087	0.016	-0.046	0.214	1.000	0.014	0.164	-0.057	0.007

Apabila dilihat dari jumlah data yang berada di kelompok yang tidak tepat, metode KHM cenderung meningkat seiring dengan bertambahnya banyak kelompok yang akan dibuat. Seperti yang dapat dilihat dari hasil di atas, nilai K=2 memiliki jumlah ketidaktepatan yang paling kecil. Hal ini berbanding lurus dengan nilai rata-rata SI yang dihasilkan.

Dari hasil tersebut, untuk pengelompokan Diktuk Bintara Polri SPN Singaraja Tahun Ajar 2018/2019 menjadi 5 kelompok dipergunakan metode KHM karena selain jumlah data yang berada dalam kelompok yang tidak tetap lebih sedikit, rata-rata nilai SI yang diperoleh juga lebih tinggi dengan komposisi seperti pada Tabel 5.

Tabel 5. Pengelompokan Anggota dengan Metode KHM

Kelompok	Rata-rata nilai centroid	Jumlah anggota	Jumlah anggota yang kurang tepat
1	84.50187	55	0
2	80.13099	20	7
3	82.24774	30	8
4	77.93042	3	0
5	82.59586	30	0

V. KESIMPULAN DAN SARAN

5.1 Kesimpulan

Berdasarkan hasil dari penelitian pengelompokan siswa Diktuk Bintara SPN Singaraja dengan menggunakan metode K-Means dan K-Harmonic Means didapatkan bahwa metode K-Harmonic Means memiliki jumlah ketidaktepatan pada pengelompokan lebih kecil dibandingkan dengan K-Means. Nilai rata-rata *silhouette indeks* yang dihasilkan oleh K-Harmonic Means juga lebih tinggi dibandingkan K-Means. Sehingga metode K-Harmonic Means sangat tepat digunakan untuk mengelompokkan siswa karena jumlah ketidaktepatan menghasilkan nilai berbanding lurus dengan nilai *silhouette indeks* yaitu jumlah ketidaktepatan lebih kecil dan nilai rata-rata *silhouette indeks* tinggi.

5.2 Saran

Berdasarkan hasil pengujian dalam data yang terbatas, kesimpulannya memang dapat dinyatakan bahwa K-Harmonic lebih baik dibandingkan K means. Terkait nilai korelasi nilai K dengan *silhouette indeks*, sampai nilai K sama dengan 9 didapatkan hubungan semakin besar nilai K, *silhouette indeks* makin kecil. Namun pada nilai K sama dengan 10, dan 11 menunjukkan hal yang sebaliknya. Pada penelitian ini pengujian hanya sampai K 11, sehingga pola yang lebih jelas belum bisa disimpulkan. Kedepan sangat penting untuk mencoba melihat korelasi nilai K yang lebih tinggi dengan *silhouette indeks*. Dengan demikian akan didapatkan formulasi yang menjelaskan korelasi yang lebih jelas.

Pada penelitian ini hanya terbatas pada uji validitas, dengan menggunakan parameter *silhouette indeks*. Kedepan sangat menarik untuk melihat kinerja algoritma K-means dan K Harmonic dari variabel lainnya, seperti bagaimana perbandingan kecepatan komputasinya.

DAFTAR PUSTAKA

- Cebeci, Z., & Yildiz, F. (2015). Comparison of K-Means and Fuzzy C-Means Algorithms on Different Cluster Structures. *Journal of Agricultural Informatics*, 6(3), 13–23.
- Gita Premashanti Trayasiwi. (2017). Penerapan Metode Klastering Dengan Algoritma K-Means Untuk Prediksi Kelulusan Mahasiswa Pada Program Studi Teknik Informatika Strata Satu. *Informaika*, 91, 399–404.
- Handoyo, R., Nasution, S. M., Studi, P., Komputer, S., Linkage, S., Coefficient, S., ... Nasution, S. M. (2014). Perbandingan Metode Clustering Menggunakan Metode Single Linkage Dan K-Means Pada Pengelompokan Dokumen. *JSM STMIK Mikroskil*, 15(2), 73–82.
- Izzah, A., & Hayatin, N. (2013). Imputasi Missing data Menggunakan Algoritma Pengelompokan Data K-Harmonic Means. *Seminar Nasional Matematika Dan Aplikasinya 2013*.
- Khan, S. S., & Ahmad, A. (2004). Cluster center initialization algorithm for K-modes clustering. *Pattern Recognition Letter*, 25(18), 1293–1302. <http://doi.org/10.1016/j.eswa.2013.07.002>
- Khan, Z., Ni, J., Fan, X., & Shi, P. (2017). An improved K-means clustering algorithm based on an adaptive initial parameter estimation procedure for image segmentation. *International Journal of Innovative Computing, Information and Control*, 13(5), 1509–1526.

- Kumar Dehariya, S. K. S. and R. C. J. (2010). Kumar Dehariya, Shailendra Kumar Shrivastava and R. C. Jain. (2010). Clustering Of Image Data Set Using K-Means And Fuzzy K-Means Algorithms. IEEE. p386-391. In *Clustering Of Image Data Set Using K-Means And Fuzzy K-Means Algorithms* (pp. 386–391). IEEE.
- Mukarromah, Martha, S., & Ihamsyah. (2015). Perbandingan Imputasi Missing Data Menggunakan Metode Mean Dan Metode Algoritma K-Means. *Buletin Ilmiah Mat. Stat. Dan Terapannya (Bimaster)*, 04(3), 305–312.
- Sani, A. (2018). Penerapan Metode K-Means Clustering Pada Perusahaan Penerapan Metode K- Means Clustering Pada Perusahaan, (August).
- Wicaksana, I. M. K., & Widiartha, I. M. (2012). Penerapan Metode Ant Colony Optimzation pada Metode K-Harmonic Means untuk Klasterisasi Data. *Jurnal Ilmu Komputer Universitas Udayana*, 5(1), 55–62.
- Zhang, B., Hsu, M., & Dayal, U. (1999). K-Harmonic Means -- A Data Clustering Algorithm,. *K-Harmonic Means - A Data Clustering Algorithm*, 1–25. Retrieved from D:\VEILLE~1\PDFBIB\ZHANG1999.PDF