

Penerapan Metode Clustering dengan Algoritma K-Means pada Pengelompokan Peminatan Mata Kuliah

Deti Karmanita

Mahasiswa Pascasarjana UPI YPTK

Billy Hendrik

Dosen Tetap Program Studi Teknik Komputer, Fakultas Ilmu Komputer,
UPI YPTK PADANG

***Abstract.** Choosing a concentration in student academic activities is not an easy thing because it depends on interests, talents and desires, therefore careful consideration is needed so that students do not make a mistake in choosing the desired concentration. This often happens when final semester students do their final assignment but it does not match their field of ability. Choosing a concentration haphazardly without careful consideration can have a negative impact on students, namely difficulty in absorbing lecture material. Therefore, a special method is needed that students can use to determine student concentration. One of the methods used is the K-Means method. The K-Means algorithm is a non-hierarchical method that initially takes a number of population components to become the initial cluster center. At this stage the cluster center is selected randomly from a set of data populations. Next, K-Means tests each component in the data population and marks the component to one of the cluster centers that has been defined depending on the minimum distance between components and each cluster. with a total of 100 data records, using cluster centers C1 70, 82.5, 85, C2 70, 75, 80 and C3 80, 85, 80 produces 6 iterations with the results of Cluster 1. Students are recommended to enter the Expert Systems Concentration. In the calculation above, there are 3 students who are included in cluster 1. Cluster 2 Students are recommended to enter the multimedia programming concentration. In the calculation above, there are 20 students included in cluster 2. Cluster 3 Students are recommended to enter the Cisci and Network Concentration. In the calculation above, there are 34 students included in cluster 3. From validation testing it is obtained: initial and final centroid of the first attribute: 5.83%, second attribute: 31.44%, third attribute: 35.89%. It is hoped to develop concentration clustering for Information Systems majors using other methods, not only the K-Means method, and determining concentration majors using variables other than academic grades, such as non-academic achievement scores which are linear with the study program. In the future, the concentration determination system will be carried out in the information systems study program.*

Keywords: Data Mining, Cluetering, K-Mens

Abstrak. Pemilihan konsentrasi dalam kegiatan akademik mahasiswa memang bukan hal yang mudah karena tergantung pada minat, bagat dan keinginan, oleh karena itu perlu pertimbangan yang matang supaya mahasiswa tidak salah dalam memilih kosentrasi yang diinginkan. Hal ini sering terjadi ketika mahasiswa semester akhir mengerjakan tugas akhir namun tidak sesuai dengan kemampuan bidang yang dimiliki. Pemilihan konsentrasi yang asal-asalan tanpa pertimbangan yang matang, menyebabkan dampak negatif pada mahasiswa, yaitu kesulitan dalam penyerapan materi-materi perkuliahan. Oleh sebab itu perlu metode khusus yang dapat digunakan mahasiswa dalam menentukan konsentrasi mahasiswa. Salah satu metode yang digunakan adalah metode K-Means. Algoritma K-Means merupakan metode non hierarki yang pada awalnya mengambil sebagian banyaknya komponen populasi untuk dijadikan pusat cluster awal. Pada tahap ini pusat cluster dipilih secara acak dari sekumpulan populasi data. Berikutnya K-Means menguji masing-masing komponen di dalam populasi data dan menandai komponen tersebut ke salah satu pusat cluster yang telah didefinisikan tergantung dari jarak minimum antar komponen dengan tiap-tiap cluster. dengan jumlah record data 100 record, menggunakan pusat cluster C1 70, 82.5, 85, C2 70, 75, 80 dan C3 80, 85, 80 menghasilkan 6 iterasi dengan hasil Cluster 1 Mahasiswa direkomendasikan masuk Konsentrasi Sistem Pakar. Dalam perhitungan di atas, ada 3 mahasiswa yang masuk dalam cluster 1. Cluster 2 Mahasiswa direkomendasikan masuk konsentrasi Pemrograman multimedia. Dalam perhitungan di atas, ada 20 mahasiswa yang masuk dalam cluster 2. Cluster 3 Mahasiswa direkomendasikan masuk Konsentrasi Cisci dan Network. Dalam perhitungan di atas, ada 34 mahasiswa yang masuk dalam cluster 3. Dari pengujian validasi diperoleh : centroid awal dan akhir atribut pertama : 5,83% , atribut kedua : 31,44% , atribut ketiga : 35,89%. Diharapkan untuk mengembangkan pengklasteran konsentrasi jurusan Sistem Informasi menggunakan metode lain tidak hanya dengan metode K-Means dan penentuan jurusan konsentrasi menggunakan variabel lain selain nilai akademik seperti nilai prestasi non akademik yang

Revised Agustus 30, 2023, Revised September 30, 2023; Accepted Oktober 23, 2023

*Deti Karmanita

linier dengan program studi. Ke depannya sistem penentuan konsentrasi yang telah dilakukan pada prodi sistem informasi.

Kata Kunci: Data Mining, Clustering, K-Means

I. PENDAHULUAN

Perguruan tinggi sebagai institusi pendidikan telah memiliki data mahasiswa dengan program studi yang berbeda sehingga dengan data tersebut dalam jumlah yang sangat besar, namun hanya sebagian kecil data tersebut dimanfaatkan (khususnya dalam penentuan peminatan jurusan).

Pemilihan konsentrasi dalam kegiatan akademik mahasiswa memang bukan hal yang mudah karena tergantung pada minat, bagat dan keinginan, oleh karena itu perlu pertimbangan yang matang supaya mahasiswa tidak salah dalam memilih konsentrasi yang diinginkan. Hal ini sering terjadi ketika mahasiswa semester akhir mengerjakan tugas akhir namun tidak sesuai dengan kemampuan bidang yang dimiliki. Pemilihan konsentrasi yang asal-asalan tanpa pertimbangan yang matang, menyebabkan dampak negatif pada mahasiswa, yaitu kesulitan dalam penyerapan materi-materi perkuliahan. Oleh sebab itu perlu metode khusus yang dapat digunakan mahasiswa dalam menentukan konsentrasi mahasiswa. Salah satu metode yang digunakan adalah metode K-Means.

Penelitian ini berbeda dari penelitian pendukung sebelumnya dalam hal manfaat dari pengelompokan jenis konsentrasi. Penelitian ini tidak hanya berfungsi memberikan penggambaran tentang konsentrasi terhadap mahasiswa namun juga memberikan penggambaran penentuan kebijakan Ketua Program studi. Sehingga dengan adanya pengelompokan data ini pihak Kaprodi juga dapat mengetahui jurusan yang paling banyak peminatnya, dari penelitian ini output yang dihasilkan adalah jenis konsentrasi yang banyak di minati oleh mahasiswa universitas dehasen bengkulu. Tujuan pada penelitian ini adalah menghasilkan sebuah data mining dalam menentukan konsentrasi jurusan bagi mahasiswa, sehingga dengan adanya sistem ini dapat membantu mahasiswa Universitas Dehasen Bengkulu dalam menentukan konsentrasi jurusan yang ada di program studi Sistem Informasi.

II. LANDASAN TIORI

a. Data Mining

Menurut Priati dan Fauzi, Ahmad (2018), definisi data mining yaitu serangkaian proses yang berguna untuk menggali nilai tambah dari suatu kumpulan data berupa pengetahuan yang tidak dapat diketahui secara manual, di mana sampel yang sama dibagi menjadi kelompok-kelompok yang disebut *cluster*.



Gambar macam-macam proses teknik data mining (Khomarudin, Agus Nur, 2016) Sebagai suatu rangkaian proses, data mining dapat dibagi menjadi beberapa tahapan . Menurut Khomarudin, Agus Nur (2016), tahapan data mining ada 6 sebagai berikut:

1. Pembersihan data (*data cleaning*), menghilangkan *noise* dan data yang tidak konsisten.
2. Integrasi data (*data integration*), menggabungkan data dari berbagai *database* ke *database* yang baru.
3. Seleksi data (*data selection*), menyeleksi semua data yang ada di *database* agar sesuai untuk dianalisis.
4. Transformasi data (*data transformation*), data diubah atau digabung ke dalam format yang sesuai.
5. Proses mining, proses yang berfungsi untuk menemukan pengetahuan berharga dan tersembunyi dari data.
6. Evaluasi pola (*pattern evaluation*), mengidentifikasi pola-pola menarik ke dalam *knowledge based* yang ditemukan.
7. Presentasi pengetahuan (*knowledge presentation*), visualisasi dan penyajian pengetahuan mengenai metode yang digunakan.

Kemajuan yang terus berlanjut dalam bidang data mining didorong oleh beberapa faktor. Menurut Intermedia, B (2020), faktor yang mendukung perlunya dilakukan data mining sebagai berikut:

1. Data telah mencapai jumlah dan ukuran yang sangat besar.
2. Telah dilakukannya proses *data warehousing*.
3. Kemampuan komputasi yang semakin terjangkau.
4. Persaingan bisnis yang semakin ketat.

Selain faktor pendukung yang diperlukan dalam proses data mining, terdapat beberapa metode data mining berdasarkan fungsi yang dilakukan atau jenis aplikasi yang dipakai yaitu sebagai berikut:

1. Klasifikasi (*supervised*), untuk segmentasi *customer*, pemodelan bisnis, analisa kartu kredit, dan yang lainnya.
2. *Clustering (unsupervised)*, untuk mengeksplorasi data.
3. *Association rules (unsupervised)*, untuk menemukan korelasi antar himpunan item.
4. *Attribute importance (unsupervised)*, untuk meningkatkan kecepatan dan akurasi dari model klasifikasi yang dibuat.

b. K-Means

K-Means adalah metode *clustering* berbasis jarak yang membagi data ke dalam sejumlah *cluster* dan algoritma ini hanya bekerja pada atribut *numeric*. Algoritma *K-Means* termasuk *partitioning clustering* yang memisahkan data ke *k* daerah bagian yang terpisah. Algoritma *K-Means* sangat terkenal karena kemudahan dan kemampuannya untuk mengelompokkan data yang besar dan data *outlier* dengan sangat cepat. Dalam algoritma *K-Means*, setiap data harus termasuk ke *cluster* tertentu dan bisa dimungkinkan bagi setiap data yang termasuk *cluster* tertentu pada suatu tahapan proses, pada tahapan berikutnya berpindah ke *cluster* lainnya.

Algoritma *K-Means* merupakan metode non hierarki yang pada awalnya mengambil sebagian banyaknya komponen populasi untuk dijadikan pusat *cluster* awal. Pada tahap ini pusat *cluster* dipilih secara acak dari sekumpulan populasi data. Berikutnya *K-Means* menguji masing-masing komponen di dalam populasi data dan menandai komponen tersebut ke salah satu pusat *cluster* yang telah didefinisikan tergantung dari jarak minimum antar komponen dengan tiap-tiap *cluster*. Posisi pusat *cluster* akan dihitung kembali sampai semua komponen data digolongkan kedalam tiap-tiap pusat *cluster* dan terakhir akan terbentuk posisi pusat *cluster* yang baru. Berikut merupakan algoritma pengelompokan data dengan metode *K-Means*.

1. Tentukan *k* sebagai jumlah *cluster* yang ingin dibentuk.
2. Bangkitkan *k centroid* (titik pusat *cluster*) awal secara random.
3. Hitung jarak setiap data ke masing-masing *centroid*.

4. Setiap data memilih *centroid* yang terdekat.
5. Tentukan posisi *centroid* yang baru dengan cara menghitung nilai rata-rata dari data-data yang terletak pada *centroid* yang sama.
6. Kembali ke langkah 3 jika posisi *centroid* baru dengan *centroid* yang lama tidak sama.

Lokasi *centroid* setiap kelompok yang diambil dari rata-rata semua nilai data pada setiap fiturnya harus dihitung kembali. Jika M menyatakan jumlah data dalam sebuah kelompok, i menyatakan fitur ke- i dalam sebuah kelompok, dan p menyatakan dimensi data (Prasetyo, 2014).

$$C_i = \frac{1}{M} \sum_{j=1}^M X_j$$

Formula tersebut dilakukan sebanyak p dimensi sehingga i mulai dari 1 sampai p . Kemudian pengukuran jarak pada data ke pusat kelompok dapat dilakukan menggunakan perhitungan jarak Euclidean menggunakan formula:

$$D(x_2, x_1) = \|x_2 - x_1\|_2 = \sqrt{\sum_{j=1}^p |x_{2j} - x_{1j}|^2}$$

D adalah jarak antara data x_2 dan x_1 , dan $| \cdot |$ adalah nilai mutlak. Metode ini mempartisi data ke dalam *cluster*/kelompok sehingga data yang memiliki karakteristik yang sama dikelompokkan ke dalam satu *cluster* yang sama dan data yang mempunyai karakteristik yang berbeda dikelompokkan ke dalam kelompok yang lain.

III. PEMBAHASAN

Data mahasiswa berdasarkan matakuliah pilihan

No	Npm	Nama	Id matkul	Mata kuliah	Sks	Nilai	Kelas	Semester
1	21030001	Herdi khrisna dwitama	MKK	Pemrograman Multimedia	2	A	A1	4
2	21030002	Setiaji ayu handayani	MKK	Sistem Pakar	2	B	A1	4
3	21030003	Endah puspa ningrum	MKK	Instalasi Jaringan Dan Internet	2	B	A1	4
4	21030004	Inda mareta indriani	MKK	Sisco Dan Network	2	A	A1	4
5	21030005	Melanisma yunita	MKK	Pemrograman Multimedia	2	A	A1	4
6	21030006	Hasna nabila	MKK	Sistem Pakar	2	B	A1	4

7	21030007	Intan mayang sari	MKK	Instalasi Jaringan Dan Internet	2	A	A1	4
8	21030008	Rafid ula penta	MKK	Sisco Dan Network	2	B	A1	4
9	21030009	Bagas fernandes	MKB	Pemrograman Multimedia	2	B	A1	4
10	21030010	Yoga pitra	MKK	Sistem Pakar	2	A	A1	4
11	21030011	Dosi martin saputra	MKK	Instalasi Jaringan Dan Internet	2	A	A1	4
12	21030012	Tiara pramuni suci	MKK	Sisco Dan Network	2	B	A1	4
13	21030013	Rupawan azzahra	MKB	Pemrograman Multimedia	2	B	A1	4
14	21030014	Qodri amril syah putra	MKK	Sistem Pakar	2	A	A1	4
15	21030015	Jeni vegas tamah	MKK	Instalasi Jaringan Dan Internet	2	B	A1	4
16	21030016	Empi desni linia junita	MKK	Sisco Dan Network	2	A	A1	4
17	21030017	Farhan zaky	MKB	Pemrograman Multimedia	2	B	A1	4
18	21030018	Dewantara	MKK	Sistem Pakar	2	A	A1	4
19	21030019	Amanda olivia firasti	MKK	Instalasi Jaringan Dan Internet	2	B	A1	4
20	21030020	Haris fadillah	MKK	Instalasi Jaringan Dan Internet	2	A	A1	4
21	21030021	Nanda novrianto	MKB	Sisco Dan Network	2	A	A2	4
22	21030022	Firda nurjanah	MKK	Pemrograman Multimedia	2	B	A2	4
23	21030023	Dimas mekar ayu fatimah	MKK	Sistem Pakar	2	A	A2	4
24	21030024	Tri ramadani	MKK	Instalasi Jaringan Dan Internet	2	A	A2	4
25	21030025	Fellco mantanno gatteris	MKB	Sisco Dan Network	2	B	A2	4
26	21030026	Raizan	MKK	Sistem Pakar	2	A	A2	4
27	21030027	Okta riansyah	MKK	Instalasi Jaringan Dan Internet	2	A	A2	4
28	21030028	Aldo	MKK	Sisco Dan	2	B	A2	4

		noprianto		Network				
29	21030029	Muhammad vinsen fahlevi	MKK	Pemrograman Multimedia	2	A	A2	4
30	21030030	Mardiansyah	MKK	Sistem Pakar	2	A	A2	4
31	21030031	Mawardi	MKK	Installasi Jaringan Dan Internet	2	B	A2	4
32	21030032	Abdul wahid	MKK	Sisco Dan Network	2	A	A2	4
33	21030033	Khairiya hayati ardis	MKK	Pemrograman Multimedia	2	B	A2	4
34	21030034	Jhon domanov	MKK	Sistem Pakar	2	A	A2	4

Simulasi Perhitungan

Simulasi perhitungan dengan Metode K_Mens dengan 10 Mahasiswa

1	21030001	Herdi khrisna dwitama	MKK	Pemrograman Multimedia	82
2	21030002	Setiaji ayu handayani	MKK	Sistem Pakar	75
3	21030003	Endah puspa ningrum	MKK	Installasi Jaringan Dan Internet	78
4	21030004	Inda mareta indriani	MKK	Sisco Dan Network	81
5	21030005	Melanisma yunita	MKK	Pemrograman Multimedia	80
6	21030030	Mardiansyah	MKK	Sistem Pakar	82
7	21030031	Mawardi	MKK	Installasi Jaringan Dan Internet	73
8	21030032	Abdul wahid	MKK	Sisco Dan Network	72
9	21030033	Khairiya hayati ardis	MKK	Pemrograman Multimedia	71
10	21030034	Jhon domanov	MKK	Sistem Pakar	74

Menentukan jumlah kelompok atau cluster Data akan dikelompokkan kedalam 3 cluster, yaitu nilai konsentrasi :

- a. Pemrograman Multimedia C1
- b. Sistem Pakar C2
- c. Instalasi Jaringan dan Internet C3
- d. Sisco dan Network C4

Menentukan pusat cluster (Centroid) Penentuan iterasi 1 dengan menentukan centroid awal, dimana centroid awal dilakukan dengan merandom dari dataset. centroid awal diambil dari data ke-1 sebagai pusat cluster 1, data ke-3 sebagai pusat cluster 2 dan data ke-5 sebagai pusat cluster 3. Berikut centroid awal setiap cluster terpapar.

Centroid Awal	A	B	C
C1	70	82.5	85
C2	70	75	80
C3	80	85	80

Menghitung iterasi pertama a.Selanjutnya proses menghitung jarak antara tiap-tiap data dengan pusat cluster (centroid). Untuk menghitung jarak tiap-tiap data terhadap setiap pusat cluster (Centroid) dapat menggunakan rumus pada persamaan (1):

Selanjutnya proses menghitung jarak antara tiap-tiap data dengan pusat cluster (centroid). Untuk menghitung jarak tiap-tiap data terhadap setiap pusat cluster (Centroid) dapat menggunakan rumus pada persamaan

$$D = \sqrt{\sum_j^n (data_i - centroid_j)^2}$$

Berikut contoh perhitungan jarak data mahasiswa pertama dengan titik pusat centroid sebagai berikut :

$$D_{1,1} = \sqrt{(70 - 70)^2 + (82.5 - 82.5)^2 + (85 - 85)^2} = 0$$

Perhitungan jarak mahasiswa pertama dengan centroid cluster kedua, seperti berikut

(3):

$$D_{1,2} = \sqrt{(70 - 70)^2 + (82.5 - 75)^2 + (85 - 80)^2} = 9.013878189$$

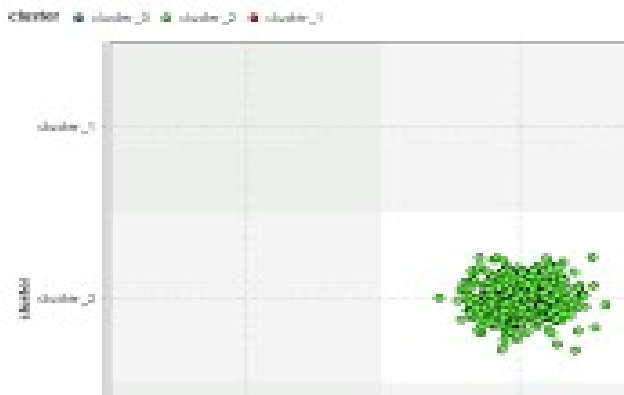
Perhitungan jarak mahasiswa pertama dengan centroid pada cluster ketiga, seperti berikut

$$D_{1,3} = \sqrt{(70 - 80)^2 + (82.5 - 85)^2 + (85 - 80)^2} = 11.45643924$$

Pada Tabel 6 berikut ini, dijelaskan hasil perhitungan jarak antar pusat cluster secara lengkap sesuai perhitungan

C1	C2	C3	JT
0	9.0138	11.4564	0
0	9.0138	11.4564	0

Gambar Clustering menggunakan rapid Miner



IV. KESIMPULAN

Hasil pengujian program terhadap data dengan jumlah record data 100 record, menggunakan pusat cluster $C1 = 70, 82.5, 85$, $C2 = 70, 75, 80$ dan $C3 = 80, 85, 80$ menghasilkan 6 iterasi dengan hasil Cluster 1 Mahasiswa direkomendasikan masuk Konsentrasi Sistem Pakar. Dalam perhitungan di atas, ada 3 mahasiswa yang masuk dalam cluster 1. Cluster 2 Mahasiswa direkomendasikan masuk konsentrasi Pemrograman multimedia. Dalam perhitungan di atas, ada 20 mahasiswa yang masuk dalam cluster 2. Cluster 3 Mahasiswa direkomendasikan masuk Konsentrasi Cisci dan Network. Dalam perhitungan di atas, ada 34 mahasiswa yang masuk dalam cluster 3. Dari pengujian validasi diperoleh : centroid awal dan akhir atribut pertama : 5,83% , atribut kedua : 31,44% , atribut ketiga : 35,89%. Diharapkan untuk mengembangkan pengklasteran konsentrasi jurusan Sistem Informasi menggunakan metode lain tidak hanya dengan metode K-Means dan penentuan jurusan konsentrasi menggunakan variabel lain selain nilai akademik seperti nilai prestasi non akademik yang linier dengan program studi. Ke depannya sistem penentuan konsentrasi yang telah dilakukan pada prodi sistem informasi.

DAFTAR PUSTAKA

- D. Dillenberger, P. Novotny, Q. Zhang, P. Jayachandran, H. Gupta, S. Hans, et al., Blockchain analytics and artificial intelligence, *IBM J. Res. Dev.* 63(2/3) (2019) 5:1–5:14.
- Darmi, Y. D., & Setiawan, A. (2016). Penerapan metode clustering k-means dalam pengelompokan penjualan produk. *Jurnal Media Infotama*, 12(2).
- Dhuhita, W. M. P. (2015). Clustering Menggunakan Metode K-Means untuk Menentukan Status Gizi Balita. *Jurnal Informatika*, 15(2), 160-174.
- IEEE standard for extensible event stream (XES) for achieving interoperability in event logs and event streams, in: *IEEE Std 1849-2016*, 2016, pp. 1–50, <http://dx.doi.org/10.1109/IEEESTD.2016.7740858>
- Intermedia, B. (2020). Data Mining : Definisi, Fungsi, Metode dan Penerapannya. *Teknologi*.
- J. Mendling, I. Weber, W.M.P.v.d. Aalst, J.v. Brocke, C. Cabanillas, F. Daniel, S. Debois, C.D. Ciccio, M. Dumas, S. Dustdar, et al., Blockchains for business process management - challenges and opportunities, *ACM Trans. Manag. Inf. Syst.* 9 (1) (2018) 4:1–4:16.
- Khomarudin, A. N. (2018). Teknik Data Mining : Algoritma K-Mean Clustering. 4-5.
- L. Moctar-M'Baba, M. Sellami, W. Gaaloul, M.F. Nanne, Blockchain logging for process mining: A systematic review, in: *55th Hawaii International Conference on System Sciences*, HICSS, Virtual Event / Maui, Hawaii, USA, January 4-7, ScholarSpace, 2022, pp. 1–10.
- N.Y. Wirawan, B.N. Yahya, H. Bae, Incorporating transaction lifecycle information in Anomali process discovery, in: *Blockchain Technology for IoT Applications*, Springer Singapore, 2021, pp. 155–172.
- Priati, & Fauzi, A. (2018). Data Mining dengan Teknik Clustering Menggunakan Algoritma K-Means pada Data Transaksi Supersore. 1.
- R. Hobeck, C. Klinkmüller, D. Bandara, I. Weber, Process mining on Anomali data: A case study of Augur, in: *19th International Conference, BPM, Rome, Italy*, Springer, 2021, pp. 306–323.
- R. Mühlberger, S. Bachhofner, C. Di Ciccio, L. García-Bañuelos, O. Pintado, Extracting event logs for process mining from data stored on the blockchain, in: *BPM International Workshops, Vienna, Austria*, in: *LNBIP*, vol.362, Springer, 2019, pp. 690–703.
- Rahajo, Budi. (2012). Modul Pemrograman Web (HTML, PHP, & MySQL). Informatika
- Selvida, D. (2019). Analisis Klasifikasi Data dengan Kombinasi Metode K-Means dan Rapid Centroid Estimation (RCE).
- Wakhidah, N. (2010). Clustering menggunakan k-means algorithm. *Jurnal Transformatika*, 8(1), 33-39.