

## Analisis Cluster dengan Data Outlier Menggunakan *Centroid Linkage* dan K-Means Clustering untuk Pengelompokan Indikator HIV/AIDS di Indonesia

Rini Silvi

Prodi Magister Statistika Terapan, Universitas Padjajaran, rinisilvi@stis.ac.id

DOI:<https://doi.org/10.15642/mantik.2018.4.1.22-31>

### Abstrak

Analisis kluster adalah salah satu metode yang digunakan untuk mengelompokkan obyek atau pengamatan yang didasarkan atas kemiripannya. Obyek yang berada dalam satu kelompok memiliki kemiripan satu sama lain. Pada penelitian ini analisis kluster yang digunakan adalah *K-means* dengan *Centroid Linkage*. *K-means* adalah salah satu metode *clustering* non-hierarki yang sederhana dan mudah diimplementasikan. Sedangkan, *Centroid Linkage* adalah metode kluster hierarki yang dapat digunakan pada data yang mengandung outlier, dimana outlier bisa membuat data yang diolah tidak mencerminkan gambaran sebenarnya. Untuk memudahkan, outlier seringkali dibuang, padahal acapkali outlier mengandung informasi penting. HIV/AIDS adalah salah satu tantangan serius terhadap kesehatan masyarakat dunia, karena HIV/AIDS merupakan penyakit menular yang menyerang sistem kekebalan tubuh sehingga penderita mengalami penurunan ketahanan tubuh secara terus-menerus yang berujung pada kematian. Data Indikator HIV/AIDS di Indonesia mengandung outlier. Penelitian ini menggunakan gap statistik untuk menentukan jumlah kluster ideal yang mengelompokkan propinsi berdasarkan indikator HIV/AIDS sedemikian hingga terbagi menjadi 7 kluster. Dari perbandingan rasio  $S_w/S_b$ , *Centroid Linkage* lebih homogen dibandingkan *K-means*. Dengan *clustering*, diharapkan pemerintah dapat mengambil kebijakan berdasarkan indikator-indikator dominan yang terdapat pada masing-masing kluster.

**Kata Kunci :** *Clustering, Centroid Linkage, K-means*

### Abstract

*Cluster analysis is a method to group data (objects) or observations based on their similarities. Objects that become members of a group have similarities among them. Cluster analyses used in this research are K-means clustering and Centroid Linkage clustering. K-means clustering, which falls under non-hierarchical cluster analysis, is a simple and easy to implement method. On the other hand, Centroid Linkage clustering, which belongs to hierarchical cluster analysis, is useful in handling outliers by preventing them skewing the cluster analysis. To keep it simple, outliers are often removed even though outliers often contain important information. HIV/AIDS is a serious challenge for global public health since HIV/AIDS is an infectious disease attacking body's immune system that in turn lowering the ability to fight infections which in the end causing death. HIV/AIDS indicators data in Indonesia contain outliers. This research uses gap statistic to define the number of clusters based on HIV/AIDS indicators that groups Indonesia provinces into 7 clusters. By comparing  $S_w/S_b$  ratio, Centroid Linkage clustering is more homogenous than K-means clustering. Using clustering, the government shall be able to create a better policy for fighting HIV/AIDS based on the dominant indicators in each cluster.*

**Keyword :** *Clustering, Centroid Linkage, K-means*

## 1. Pendahuluan

HIV/AIDS merupakan penyakit menular yang disebabkan oleh infeksi *Human Immunodeficiency Virus* yang menyerang sistem kekebalan tubuh sehingga tidak dapat melawan infeksi dan berujung pada kematian. Penyakit ini merupakan masalah dan tantangan serius terhadap kesehatan masyarakat di dunia. Di Indonesia, sebagian besar infeksi baru diperkirakan terjadi pada beberapa sub-populasi berisiko tinggi dengan prevalensi > 5% [1]. Risiko penularan HIV/AIDS tidak hanya terbatas pada sub-populasi yang berperilaku risiko tinggi, tetapi juga dapat menular pada pasangan atau bahkan anaknya.

Pada tahun 2016, kasus baru infeksi HIV meningkat 33,4% dibandingkan tahun sebelumnya. Terdapat 69,3% kasus baru infeksi HIV pada kelompok umur 25-49, sementara 63,3% penderita adalah laki-laki. Rasio HIV/AIDS antara laki-laki dan perempuan tercatat pada kisaran 2:1 [2]. Pada tahun 2016, jumlah kasus baru infeksi HIV terbanyak adalah Jawa Timur, kemudian diikuti oleh Jakarta, Jawa Barat, Jawa Tengah, dan Papua. Akan tetapi, persentase kasus baru infeksi HIV terbesar adalah Papua, jika dibagi dengan jumlah penduduknya. Sementara, Gorontalo menempati urutan terakhir propinsi dengan jumlah kasus baru infeksi HIV terbanyak di Indonesia.

Perbedaan jumlah kasus baru HIV tampak begitu nyata di beberapa propinsi. Contohnya DKI Jakarta dengan tingkat prevalensi 58,56 padahal dengan tingkat kemiskinan paling kecil yaitu sebesar 3,73. Meskipun tingkat pendidikan di Gorontalo termasuk rendah, akan tetapi tingkat prevalensi kasus baru HIV hanya sekitar 0,6. Tingkat prevalensi disini dihitung berdasarkan jumlah kasus baru infeksi HIV suatu daerah dibandingkan setiap seratus ribu penduduk.

Terdapat beberapa indikator yang mempengaruhi HIV/AIDS yaitu tingkat penggunaan kontrasepsi (kondom), jumlah dokter/tenaga medis, proporsi muslim, tingkat kesuburan remaja, dan rata-rata lama sekolah [3]. Sementara menurut [4], HIV dipengaruhi oleh indikator seperti kurangnya pendidikan, kemiskinan, seks bebas, dan kehidupan malam meskipun pada penelitiannya lebih difokuskan kepada perempuan.

Ada beberapa indikator yang dapat mempengaruhi prevalensi HIV yang digunakan disini. Pada indikator tersebut teridentifikasi

beberapa data outlier. Hal ini bisa mempengaruhi hasil kesimpulan dari penelitian. Outlier dapat menghasilkan output yang tidak sesuai dengan gambaran yang sebenarnya, termasuk dalam hal pengklasteran indikator prevalensi HIV. Maka dari itu, penelitian ini perlu menggunakan metode klaster yang dapat menangani pengaruh keberadaan outlier tersebut. Metode yang dimaksud adalah K-Means Clustering.

Tujuan dari penelitian ini adalah untuk mengumpulkan propinsi ke dalam suatu kelompok sedemikian hingga dapat dibedakan menurut karakteristik indikator dari prevalensi HIV. Dengan adanya kelompok-kelompok berdasarkan indikator-indikator tersebut, diharapkan pemerintah dapat membuat kebijakan yang tepat untuk pencegahan dini menyebarnya HIV/AIDS.

## 2. Tinjauan Pustaka

### 2.1 Analisis Klaster

Analisis klaster adalah suatu teknik statistik yang bertujuan untuk mengelompokkan obyek ke dalam suatu kelompok sedemikian sehingga obyek yang berada dalam satu kelompok akan memiliki kesamaan yang tinggi dibandingkan dengan obyek yang berada di kelompok lain [5]. Dengan kata lain tujuan dari analisis cluster adalah pengklasifikasian obyek-obyek berdasarkan similaritas diantaranya dan menghim-pun data menjadi beberapa kelompok [16]. Ada dua metode dalam analisis klaster yaitu metode hierarki dan metode non hierarki. Menurut [6], metode non hierarki umumnya digunakan jika jumlah satuan pengamatan besar dan jumlah klaster tidak ditentukan sebelumnya. Salah satu metode non hierarki adalah metode K-means. Ini adalah metode non hirarki yang paling banyak digunakan. Algoritma K-means mudah diimplementasikan dan juga mudah diadaptasi sehingga menjadikannya lebih populer dalam hal pengelompokan. Pada teknik K-means, biasanya peneliti sudah terlebih dahulu menentukan banyaknya klaster yang akan dibentuk.

Metode hierarki merupakan metode pengelompokan yang terstruktur dan bertahap berdasarkan kemiripan sifat antar obyek. Kemiripan sifat tersebut dapat ditentukan dari kedekatan jarak Euclidean atau jarak Mahalanobis.

Jarak Euclidean digunakan jika tidak terjadi korelasi. Jarak Euclidean dirumuskan sebagai berikut:

$$d(y, x) = \sqrt{\sum_{k=1}^l (y_k - x_k)^2}; l = 1, 2, 3, \dots, n$$

dimana:

$d(y, x)$  : kuadrat jarak Euclid antar obyek  $y$  dengan obyek pada  $x$ .

$y_k$  : nilai dari obyek  $y$  pada variabel ke- $k$

$x_k$  : nilai dari obyek  $x$  pada variabel ke- $k$

Jarak Mahalanobis digunakan jika data terjadi korelasi. Jarak Mahalanobis antara dua sampel  $X$  dan  $Y$  dari suatu variabel acak didefinisikan sebagai berikut [5] :

$$d_{Mahalanobis}(y, x) = \sqrt{(y - x)^T \Sigma^{-1} (y - x)}$$

Dengan  $\Sigma$  adalah suatu matriks varians kovarians.

Dalam metode kluster hierarki terdapat beberapa metode penghitungan jarak yang dapat digunakan, antara lain metode pautan tunggal (*single linkage*), metode pautan lengkap (*complete linkage*), metode pautan rata-rata (*average linkage*), metode Ward, dan metode *Centroid Linkage*.

## 2.2 Multikolinearitas

Multikolinearitas adalah adanya hubungan linear yang sempurna atau pasti di antara beberapa atau semua variabel [7]. Multikolinearitas berkenaan dengan terdapatnya lebih dari satu hubungan linear pasti. Untuk mengetahui adanya multikolinearitas salah satunya adalah dengan menghitung nilai *Variance Inflation Factor* (VIF) dengan rumus:

$$VIF_j = \frac{1}{1 - R_j^2}$$

Menurut [7], terjadinya multikolinearitas apabila nilai ( $VIF_j$ )  $\geq 10$ . Jika terindikasi terjadi multikolinearitas maka harus dilakukan tindakan perbaikan multikolinearitas.

## 3. Metode Penelitian

### 3.1 Sumber Data dan Variabel Penelitian

Penelitian ini menggunakan data sekunder tahun 2016 yang diperoleh dari Badan Pusat Statistik Republik Indonesia (BPS RI) kemudian diolah dengan *software R* versi 3.4.3. Data prevalensi kasus baru infeksi HIV digunakan sebagai data pembanding dalam

pengelompokkan kluster. Variabel indikator yang digunakan menjadi faktor risiko prevalensi HIV yaitu persentase penduduk miskin ( $X_1$ ) [8], tingkat pengangguran terbuka ( $X_2$ ) [9], jumlah puskesmas ( $X_3$ ) [10] dan rasio penduduk 15 tahun ke atas yang masih berpendidikan rendah ( $X_4$ ) [11], dan persentase pasangan usia subur yang memakai alat kontrasepsi berupa kondom ( $X_5$ ) [2]. Pendidikan rendah diasumsikan untuk yang berumur 15 tahun keatas tapi tidak punya ijazah, atau hanya sampai tamatan SD dan SMP. Alat kontrasepsi yang digunakan difokuskan pada kondom karena alat kontrasepsi yang lainnya hanya digunakan untuk mencegah kehamilan tetapi tidak bisa mengurangi penularan HIV/AIDS. Unit observasi pada penelitian ini adalah seluruh propinsi di Indonesia.

## 3.2 Tahapan Penelitian

### 3.2.1 Standardisasi Data

Proses standardisasi dilakukan apabila di antara variabel-variabel yang diteliti terdapat perbedaan ukuran satuan yang besar. Perbedaan satuan yang mencolok dapat mengakibatkan perhitungan pada analisis kluster menjadi tidak valid. Oleh karena itu, perlu dilakukan proses standardisasi dengan melakukan transformasi pada data asli sebelum dianalisis lebih lanjut.

### 3.2.2 Deteksi Outlier dan Multikolinearitas

Analisis kluster pada hakekatnya adalah teknik algoritma, bukan alat inferensi statistik. Oleh sebab itu persyaratan seperti distribusi data yang harus normal (di analisis statistik lainnya) ataupun hubungan linier antar variabel tidak menjadi syarat dalam analisis kluster. Namun demikian, karena data yang diolah dalam analisis kluster biasanya hanya sebagian kecil dari populasi, maka agar hasilnya bisa digeneralisasi, data yang diolah sebaiknya mencerminkan gambaran umum atau bersifat representatif. Oleh sebab itu, *outliers* tetap harus dihilangkan dari sampel agar hasilnya tidak bias.

Deteksi outlier digunakan untuk mencari data yang berbeda dengan mayoritas data yang lain. Walaupun memiliki perilaku yang berbeda dengan mayoritas data yang lain dan sering dianggap *noise*, tetapi outlier sering kali mengandung informasi yang sangat berguna. Tidak semua data yang mengandung outlier bisa

ditransformasi karena kasus data yang berbeda-beda. Akan tetapi, dengan menggunakan metode *Centroid Linkage*, outlier tidak berpengaruh secara signifikan.

Selain itu, data yang digunakan seharusnya tidak berkorelasi, dengan kata lain sebaiknya tidak ada multikolinieritas. Alasannya adalah di dalam analisis kluster setiap variabel diberi bobot yang sama dalam perhitungan jarak. Manakala beberapa variabel saling berkorelasi, korelasi tersebut akan menyebabkan pembobotan yang tidak berimbang sehingga akan mempengaruhi hasil analisis [12].

### 3.2.3 Penentuan Jumlah Kluster Optimum

Penentuan jumlah *cluster* optimum dilakukan dengan menggunakan gap statistik pada  $R$ . Gap statistik bertujuan untuk menentukan jumlah kluster lebih konstan dibandingkan pengukuran lainnya. Jarak obyek berpasangan dalam kluster dihitung dengan rumus:

$$D_r = \sum_{i, i' \in C_r} d_{ii'}$$

Dimana  $d$  adalah jarak euclidean kuadrat. Kemudian hitung jumlah kuadrat di dalam kluster menggunakan rumus:

$$W_k = \sum_{r=1}^k \frac{1}{2n_r} D_r$$

Nilai gap didapatkan dengan mengestimasi jumlah kluster optimum pendekatan standardisasi  $W_k$ :

$$Gap_n(k) = E_n^* \{ \log(W_k) \} - \log(W_k)$$

Dimana  $E_n^*$  adalah nilai ekspektasi dari distribusi jumlah sampel.

Kriteria banyak *cluster* optimum diberikan oleh nilai gap statistik ( $k$ ) yang paling tinggi, atau yang pertama kali mengindikasikan kenaikan gap yang minimum jika gap selalu naik [13]. Setelah penentuan kluster optimum, maka akan dibandingkan pengklasteran dengan menggunakan dua metode, yaitu metode *Centroid Linkage* dan metode K-means.

### 3.2.4 Metode Kluster K-means

Metode K-means adalah metode non hierarki yang paling banyak digunakan dalam pengklasteran. Algoritma K-means mudah diimplementasikan. Pada metode ini, peneliti menentukan sendiri jumlah kluster yang akan dibentuk. Peneliti mengelompokkan entitas ke dalam k-kelompok, biasanya dilakukan secara

acak. Pada masing-masing kelompok dihitung rata-ratanya. Hitung jarak setiap entitas terhadap pusat masing-masing kelompok (rata-rata kelompoknya).

Masing-masing objek dialokasikan ke kluster terdekat dengan pusatnya. *Update* keanggotaan setiap entitas berdasarkan jarak terdekat dengan pusat kelompok dan ditentukan kembali pusat kluster yang baru. Proses pengalokasian obyek kembali dilakukan. Suatu obyek dapat berpindah ke kluster lain bila obyek tersebut lebih dekat ke pusat kluster tersebut. Proses ini dilakukan secara berulang sampai tidak ada lagi entitas yang berpindah kelompok.

### 3.2.5 Metode Kluster Centroid Linkage

*Centroid Linkage* adalah rata-rata semua obyek dalam kluster. Jarak antara dua kluster adalah jarak antar *centroid* kluster tersebut. Kluster *centroid* adalah nilai tengah observasi pada variabel dalam suatu set variabel cluster. Dengan metode ini, setiap terjadi kluster baru segera terjadi perhitungan ulang *centroid* sampai terbentuk kluster yang tetap [5].

Keuntungan dari metode ini adalah outlier tidak berpengaruh secara signifikan, jika dibandingkan dengan metode lain. Jarak antara dua kluster didefinisikan sebagai berikut:

$$d_{(UV)W} = d(\bar{x}_1, \bar{x}_2)$$

Centorid kluster baru yang terbentuk didapat dengan rumus:

$$\bar{x} = \frac{N_1 \bar{x}_1 + N_2 \bar{x}_2}{N_1 + N_2}$$

Dimana:

$$N_1 = N_2 : \text{banyaknya obyek.}$$

*Centroid* adalah rata-rata dari semua anggota dalam kluster tersebut. Pada saat obyek digabungkan maka *centroid* baru dihitung, sehingga setiap kali ada penambahan anggota, *centroid* akan berubah pula [14].

### 3.2.6 Langkah-langkah K-means Clustering

Langkah-langkah K-means:

1. Menentukan k sebagai jumlah kluster yang diinginkan
2. Mengalokasikan data ke dalam cluster secara acak
3. Menentukan pusat kluster dari data yang ada pada masing-masing kluster dengan persamaan:

$$C_{kj} = \frac{x_{1j} + x_{2j} + \dots + x_{nj}}{n}$$

dimana:

$C_{kj}$  : pusat *cluster* ke-k pada variabel ke j

(j=1, 2, ..., p)

n : banyak data pada *cluster* ke-k

4. Menentukan jarak setiap obyek dengan setiap *centroid* dengan perhitungan jarak setiap obyek dengan setiap *centroid* menggunakan jarak *Euclidean*
5. Menghitung fungsi obyektif dengan formula:
 
$$l = \sum_{i=1}^n \sum_{j=1}^k a_{ij} d(x_i, C_{kj})^2$$
6. Mengalokasikan masing-masing data ke *centroid*/rata-rata terdekat.
7. Mengulangi kembali langkah 3-6 sampai tidak ada lagi perpindahan obyek atau tidak ada perubahan pada fungsi obyektifnya.

### 3.2.7 Langkah-langkah Centroid Linkage Clustering

Langkah-langkah *Centroid Linkage Clustering*:

1. Membuat k kluster. Masing-masing individu atau unit observasi jadi kelompok. Kemudian dibuat matrik jaraknya (dari i ke kelompoknya), dengan rumus:

$$D = \{d_{ik}\}$$

2. Mencari jarak terkecil dari pasangan kluster, yaitu  $d_{uv}$  (jarak kluster u dan kluster v)
3. Menggabungkan kluster u dan kluster v. kemudian update matriks jaraknya.
4. Ulangi langkah 2 & 3 sebanyak N-1 kali. Catat nilai jarak untuk setiap terjadi penggabungan kluster
5. Tentukan nilai *cut off* untuk menentukan kluster terbentuk. Lakukan dengan membuat dendogram, dan tentukan *cut off* jumlah kluster
6. Pemberian nama kluster berdasarkan *profiling*, yaitu melihat karakteristik kluster terbentuk secara rata-rata.

### 3.2.8 Penentuan Metode Terbaik

Tahap Evaluasi dapat dilakukan dengan melakukan analisis kluster dengan ukuran jarak atau metode kluster yang berbeda kemudian dibandingkan hasilnya [7]. Pemilihan metode yang menghasilkan kualitas pengelompokan terbaik dilakukan dengan memperhatikan nilai rasio rata-rata simpangan baku dalam kluster terhadap simpangan baku antar kluster [15]. Rata-rata simpangan baku di dalam kluster ( $S_w$ ) dinyatakan dengan:

$$S_w = \frac{1}{c} \sum_{k=1}^c S_k$$

Simpangan baku antar kluster ( $S_b$ ) dinyatakan sebagai:

$$S_b = \left[ \frac{1}{c-1} \sum_{k=1}^c (\bar{X}_k - \bar{X})^2 \right]^{\frac{1}{2}}$$

Dimana c adalah jumlah kluster,  $S_k$  merupakan simpangan baku di dalam kluster ke-k.  $\bar{X}_k$  sebagai rata-rata kluster ke-k dan  $\bar{X}$  adalah rata-rata dari semua kluster. Semakin kecil nilai  $S_w$  dan semakin besar nilai  $S_b$  maka metode tersebut memiliki kinerja yang baik, artinya memiliki homogenitas yang tinggi. Metode yang dipilih adalah yang memberikan nilai rasio  $S_w/S_b$  terkecil.

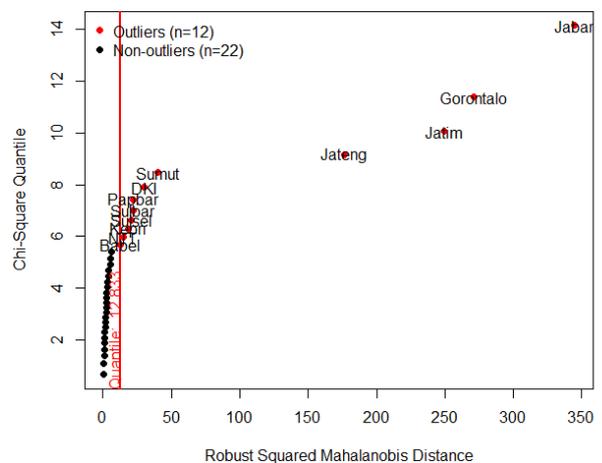
## 4. Hasil dan Pembahasan

### 4.1 Jumlah Kluster Optimum

Sebelum melakukan analisis lebih lanjut terlebih dahulu dilakukan standarisasi data, karena terdapat perbedaan satuan yang mencolok sehingga dapat mengakibatkan perhitungan pada analisis kluster menjadi tidak valid. Transformasi dilakukan menggunakan R GUI 3.4.3.

Pemeriksaan awal yang dilakukan adalah melihat apakah variabel yang digunakan terdapat outlier atau tidak. Pada data ini terdapat beberapa provinsi yang outlier, diantaranya yaitu Jawa Barat, Gorontalo, dan Jawa Timur, Jawa Tengah, serta beberapa daerah lainnya. Pada Gambar 1, terdapat 12 titik outlier. Untuk menangani masalah outlier, tanpa melakukan perubahan pada data, perlu dibandingkan metode terbaik dari beberapa metode yang ada.

Chi-Square Q-Q Plot



Gambar 1. Pendeteksian Outlier

Mendeteksi multikolinearitas dengan menggunakan VIF dengan hasil masing-masing





Klaster kelima, adalah daerah dengan proporsi pasangan usia subur (PUS) dengan kontrasepsi kondom relatif kecil, akan tetapi jumlah puskesmas yang terdapat di daerah ini termasuk sedikit. Klaster ini adalah daerah Nusa Tenggara Timur.

Klaster keenam, adalah daerah dengan proporsi pasangan usia subur (PUS) dengan kontrasepsi kondom paling banyak, sehingga klaster ini termasuk daerah dengan jumlah HIV/AIDS terkecil, yaitu daerah Gorontalo.

Klaster ketujuh, adalah daerah dengan persentase penduduk miskin paling besar dan pendidikan paling rendah. Yaitu daerah Papua.

## 5. Kesimpulan

Hasil penelitian ini memberikan kesimpulan bahwa untuk data yang memiliki outlier, metode pengklasteran menggunakan *Centroid Linkage* lebih memberikan hasil yang sesuai dengan keadaan dibandingkan dengan metode K-means. Metode K-means lebih heterogen dalam hal ini. Dengan metode *Centroid Linkage*, outlier tidak mempengaruhi klaster analisis dan tidak mengubah hasil dari interpretasi data.

Dari 34 propinsi yang ada di Indonesia, terdapat 7 klaster berdasarkan indikator yang menyebabkan terjadinya HIV/AIDS.

Klaster 1: Aceh, Sumatera Utara, Sumatera Barat, Riau, Jambi, Sumatera Selatan, Bengkulu, Lampung, Bangka Belitung, DIY, Banten, Bali, NTB, Kalimantan Barat, Kalimantan Tengah, Kalimantan Selatan, Kalimantan Timur, Kalimantan Utara, Sulawesi Utara, Sulawesi Tengah, Sulawesi Selatan, Sulawesi Tenggara, Sulawesi Barat, Maluku, Maluku Utara, Papua Barat.

Klaster 2: Kepulauan Riau, dan DKI Jakarta.

Klaster 3: Jawa Barat.

Klaster 4: Jawa Tengah dan Jawa Timur.

Klaster 5: Nusa Tenggara Timur

Klaster 6: Gorontalo

Klaster 7: Papua

## Referensi

- [1] Kementerian Kesehatan RI, Survei Terpadu Biologis dan Perilaku, Jakarta: Kementerian Kesehatan, (2011)
- [2] Kementerian Kesehatan RI, Profil Kesehatan Indonesia Tahun 2016, Jakarta: Kementerian Kesehatan RI, (2017)
- [3] Mondal, M., & Shitan, M., Factors Affecting The HIV/AIDS Epidemic: An Ecological Analysis Of Global Data, *African Health Sciences*, 13(2) pp 301–310, (2013)
- [4] Singh, R. K., & Patra, S., What Factors Are Responsible For Higher Prevalence Of HIV Infection Among Urban Women Than Rural Women In Tanzania?, *Ethiopian Journal of Health Sciences*, 25(4), pp 321–328, (2015)
- [5] Rahmawati, L., Analisis Kelompok Dengan Menggunakan Metode Hierarki Untuk Pengelompokan Kabupaten/Kota Di Jawa Timur Berdasar Indikator Kesehatan, *Jurnal Matematika Vol.1 No.2 Universitas Negeri Malang*, (2012)
- [6] Anderberg, M., *Cluster Analysis For Applications*, Academic Press.Inc., (1973)
- [7] Ningrat, D.R., Analisis Cluster Dengan Algoritma K-Means dan Fuzzy C-Means Clustering Untuk Pengelompokan Data Obligasi Korporasi, *Jurnal Gaussian Vol.5 No.4 Universitas Diponegoro*, 2016.
- [8] Badan Pusat Statistik (BPS), *Indikator Pembangunan Berkelanjutan 2017: Jumlah Penduduk Miskin*, (2017)
- [9] Badan Pusat Statistik (BPS), *Statistik Indonesia 2017: Tingkat Pengangguran Terbuka*, (2017)
- [10] Badan Pusat Statistik (BPS), *Statistik Indonesia 2017: Jumlah Puskesmas*, (2017)
- [11] Badan Pusat Statistik (BPS), *Statistik Kesejahteraan Rakyat 2017: Status Pendidikan Tertinggi*, (2017)
- [12] Puspitasari, M., *Pengelompokan Kabupaten / Kota Berdasarkan Faktor-Faktor Yang Mempengaruhi Kemiskinan Di Jawa Tengah Menggunakan Metode Ward Dan Average Linkage*, *Jurnal Matematika Vol. 5 No. 6 Universitas Negeri Yogyakarta*, (2016)
- [13] Tibshirani, R., Walther, G., & Hastie, T., Estimating The Number Of Clusters In A Data Set Via The Gap Statistic, *Journal of Royal Statistical Society Vol. 63 Issue 2.*, (2001)
- [14] Laeli, S., Analisis Cluster Dengan Average Linkage Method Dan Ward's Method Untuk Data Responden Nasabah Asuransi Jiwa Unit Link, S1 Thesis, Universitas Negeri Yogyakarta, Indonesia, (2014)
- [15] Purnamasari, S.B., Pemilihan Cluster Optimum Pada Fuzzy C-Means (Studi Kasus: Pengelompokan Kabupaten/Kota di

Jawa Tengah berdasarkan Indikator Indeks Pembangunan Manusia), Jurnal Gaussian Vol.3 No.3 Universitas Diponegoro, (2014)

[16] Lailiyah, S. dan Hafiyusholeh, M., Perbandingan antara Metode K-Means Clustering dengan Gath-Geva Clustering, Jurnal Matematika MANTIK, 1(2), Mei 2016. pp. 26-37

## Lampiran

Tabel 3. Prevalensi Penderita HIV per100.000 Penduduk Menurut Propinsi Tahun 2016

Provinsi	Rasio Penderita HIV*100.000	Provinsi	Rasio Penderita HIV*100.000
Aceh	1.37	Nusa Tenggara Barat	3.59
Sumatera Utara	13.41	Nusa Tenggara Timur	9.36
Sumatera Barat	7.53	Kalimantan Barat	10.80
Riau	12.64	Kalimantan Tengah	5.53
Jambi	6.22	Kalimantan Selatan	11.19
Sumatera Selatan	4.24	Kalimantan Timur	23.22
Bengkulu	6.04	Kalimantan Utara	24.46
Lampung	4.64	Sulawesi Utara	16.78
Kep. Bangka Belitung	9.63	Sulawesi Tengah	5.37
Kepulauan Riau	51.13	Sulawesi Selatan	11.54
DKI Jakarta	58.56	Sulawesi Tenggara	5.25
Jawa Barat	11.54	Gorontalo	0.61
Jawa tengah	11.85	Sulawesi Barat	1.68
DI Yogyakarta	19.78	Maluku	36.20
Jawa Timur	16.67	Maluku Utara	10.12
Banten	8.95	Papua Barat	59.32
Bali	56.36	Papua	120.53