

BAB III

METODOLOGI PENELITIAN

3.1. Objek Penelitian

Objek penelitian pada penelitian ini adalah penyebaran Covid-19 di DKI Jakarta. Dimana peneliti memakai informasi mengenai kasus Covid-19 di tahun 2020 – 2022 dengan perincian yaitu bulan Maret – Desember 2020, Januari – Desember 2021, Januari – November 2022, selain itu penelitian yang akan dilakukan memiliki tujuan yaitu mengelompokkan penyebaran Covid-19 di DKI Jakarta. Hasil penelitian ini diharapkan akan bisa digunakan sebagai sebuah masukan untuk menentukan kebijakan terhadap masing-masing daerah di DKI Jakarta guna menekan angka penyebaran Covid-19 dengan data terbaru yaitu tahun 2020 sampai 2022. Data tersebut dikategorikan dengan melihat mean terdekat, sehingga hasil penelitian adalah kluster dari masing-masing daerah yang terkena dampak Covid-19.

3.2. Metode Penelitian

Metode *k-means* klustering dan *fuzzy c-means* digunakan untuk membentuk klusterisasi penyebaran Covid-19 di DKI Jakarta. Metode *k-means* dipilih karena mudah diterapkan dan dijalankan, relatif cepat, mudah diadaptasi serta banyak dipraktikkan dalam tugas data mining [27]. Selain itu, nilai akhir yang diperoleh dari metode *k-means* dapat memaksimalkan kemiripan data pada satu kluster dan meminimalkan kemiripan data antar kluster [14]. Sedangkan, metode *fuzzy c-means* dipilih karena memiliki kelebihan yang terletak pada penempatan pusat kluster yang lebih tepat dibandingkan dengan metode kluster lain [28].

3.2.1. *K-Means Klustering*

Termasuk bagian dari metode klustering *non-hirarki* yang dipakai untuk pengklasteran data ke satu atau beberapa kelompok. Penerapan metode ini yaitu ketika data dengan ciri yang serupa akan diklasterkan dalam sebuah kelompok,

sedangkan data yang karakteristiknya berbeda dikelompokkan dengan kelompok yang berbeda.

3.2.2. Fuzzy C-Means

Algoritma *Fuzzy C-Means* merupakan teknik pengelompokan yang terawasi, karena pada algoritma ini jumlah kluster yang akan dibentuk perlu diketahui terlebih dahulu. Konsep dasar algoritma *fuzzy c-means* adalah menentukan pusat kelompok yang akan menandai lokasi rata-rata untuk setiap tiap-tiap kluster. Tujuan dari algoritma *fuzzy c-means* untuk mendapatkan pusat klaster yang nantinya akan digunakan untuk mengetahui data yang masuk ke dalam sebuah kluster [29].

3.2.3. Perbandingan Metode

Metode yang digunakan dalam penelitian ini adalah *K-means Klustering* dan *fuzzy c-means*. Kedua metode ini dipakai untuk melakukan pengklasteran. Kedua metode tersebut memiliki karakteristik yang berbeda dalam penerapannya.

Tabel 3.1 Perbandingan Metode

Parameter	<i>K-Means Klustering</i>	<i>Fuzzy C-Means</i>
Fase	<p>Tahapan <i>K-Means</i> [30]:</p> <ol style="list-style-type: none"> 1. Tentukan k sebagai jumlah kluster yang dibentuk. 2. Tentukan k centroid (titik pusat kluster) awal secara random. 3. Hitung jarak setiap objek ke masing-masing centroid dari masing-masing kluster. 4. Alokasi masing-masing objek ke dalam centroid yang paling dekat. 5. Lakukan iterasi, kemudian tentukan posisi centroid baru. 6. Ulangi langkah 3 jika posisi centroid baru tidak sama. 	<p>Tahapan <i>Fuzzy C-Means</i> [31]:</p> <ol style="list-style-type: none"> 1. Inisialisasi K (Jumlah kluster), $f(x)$ (fungsi objektif) awal, centroid awal (K data sebagai centroid awal), dan tetapkan ambang batas (threshold), pembobot (w) dan maksimum iterasi. 2. Hitung nilai centroid tiap-tiap kluster. 3. Hitung nilai derajat keanggotaan setiap data ke setiap Kluster. 4. Hitung $f(x)$ dengan metode perhitungan jarak Euclidean. 5. Perbaiki derajat keanggotaan setiap data pada setiap kluster. 6. Ulangi langkah 2 dan 5 hingga konvergen tercapai, yaitu apabila derajat keanggotaan \leq ambang batas, atau perubahan $f(x) \leq$ ambang batas, atau apabila

		perubahan centroid \leq ambang batas, atau telah mencapai iterasi maksimum yang ditentukan di awal.
Kelebihan	Kelebihan <i>K-Means</i> [30]: <ol style="list-style-type: none"> 1. Metode yang sangat simple dan fleksibel. Artinya perhitungan komputasinya tidak telalu rumit dan dapat diimplementasikan pada segala bidang. 2. Metode ini sangat mudah dipahami, terutama dalam implementasi data yang sangat besar serta dapat mengurangi kompleksitas data yang dimiliki. 	Kelebihan <i>Fuzzy C-Means</i> [32]: <ol style="list-style-type: none"> 1. Metode ini bersifat <i>unsupervised</i>. 2. Dapat mencapai pusat <i>klaster</i> yang <i>konvergen</i>. 3. Dalam kondisi tertentu FCM merupakan model klustering yang mempunyai ketangguhan jika dilihat dari nilai fungsi obyektifnya, jumlah iterasinya dan waktu yang diselesaikan.
Kekurangan	Kekurangan <i>K-Means</i> [30]: <ol style="list-style-type: none"> 1. Memerlukan angka yang tepat dalam menentukan jumlah klaster sebanyak <i>k</i>. Karena terkadang pusat klaster awal dapat berubah sehingga kejadian ini bisa mengakibatkan pengelompokkan data menjadi tidak stabil. 2. Metode ini tidak bisa maksimal dalam menentukan atau menginisialisasi nilai centroid awalnya, karena pada pengelompokkan data dengan metode <i>K-Means</i> sangat bergantung pada nilai centroid-nya. 3. Ditemukannya beberapa model klaster yang berbeda. 4. Harus melakukan pemilihan jumlah klaster yang tepat. 5. Kegagalan untuk converge 6. Input dari <i>K-Means</i> tergantung pada nilai pusat yang dipilih pada <i>Klustering</i>. 	Kekurangan <i>Fuzzy C-Means</i> [32]: <ol style="list-style-type: none"> 1. Mudah terjebak dalam <i>local optima</i>. 2. Sensitif terhadap pusat klaster awal.
Hasil Penerapan	Output dari <i>K-Means</i> [33]: <ol style="list-style-type: none"> 1. Klaster: vector yang berisikan lokasi klaster tiap objek. 2. <i>Centers</i>: matriks yang berisikan centroid/rata-rata nilai tiap klaster. 3. <i>Withinss</i>: vektor yang berisikan simpangan tiap klaster yang terbentuk. 4. Jumlah objek pada tiap <i>klaster</i>. 	Berbentuk deret dari pusat klaster dan keanggotaan terhadap setiap titik data. Hasil akan dipakai dalam membangun <i>fuzzy interface system</i> [34].

Perbandingan metode penelitian yang umum digunakan pada proses data mining, yaitu metode KDD, CRISP-DM dan SEMMA.

Tabel 3.2 Perbandingan Metode *Data Mining* [35]

	KDD	CRISP-DM	SEMMA
Tahapan	<ol style="list-style-type: none"> 1. <i>Data selection</i> 2. <i>Pre-processing</i> atau <i>cleaning</i> 3. <i>Transformation</i> 4. <i>Data mining</i> 5. <i>Interpretation</i> atau <i>Evaluation</i> 	<ol style="list-style-type: none"> 1. <i>Business understanding</i> 2. <i>Data understanding</i> 3. <i>Data preparation</i> 4. <i>Modeling</i> 5. <i>Evaluation</i> 6. <i>Deployment</i> 	<ol style="list-style-type: none"> 1. <i>Sample</i> 2. <i>Explore</i> 3. <i>Modify</i> 4. <i>Model</i> 5. <i>Assessment</i>

Dari tabel 3.2 mengenai perbandingan metode penelitian *data mining* tersebut pada tabel diatas, peneliti menetapkan penggunaan metode CRISP-DM (*Cross-Industry Standard Process for Data Mining*) pada penelitian ini dikarenakan CRISP-DM sebagai salah satu standard untuk menghasilkan *data driven decision making* yang berkualitas berdasarkan hasil *pooling* datascience-pm, yang dimana CRISP-DM digunakan 2 sampai dengan 3 kali lebih banyak dari 5 teratas standard yang paling banyak digunakan. Metode CRISP-DM dipilih karena model proses standar terbuka yang menggambarkan pendekatan umum yang digunakan untuk *data mining*. Selain itu, CRISP-DM bertujuan untuk menyediakan proses *data mining* yang andal dan berulang. Hal tersebut mempunyai kerangka kerja yang seragam dengan pedoman dan fleksibel untuk beradaptasi dengan masalah bisnis dan kumpulan data yang berbeda [36].

3.3. Teknik Pengumpulan Data

Penelitian ini memakai pendekatan kuantitatif. Tahapan yang digunakan dalam pengumpulan data adalah observasi untuk memperoleh data yang diperlukan berupa data melalui website Corona Jakarta dan kemudian dilakukan *scrapping*.

3.4. Teknik Pengambilan Sampel

Teknik *Purposive Sampling* digunakan dalam penelitian ini untuk pengambilan sampel. Teknik *purposive sampling* merupakan metode jenis *sampling*

non-random yang dilakukan dengan mempertimbangkan hal tertentu dalam pengambilan sampel untuk tujuan tertentu [37]. Sampel yang digunakan adalah data Covid-19 di DKI Jakarta pada tahun 2020-2022 yang diambil melalui laman corona.jakarta.go.id.

3.5. Variabel Penelitian

Variabel penelitian menggunakan data tahun 2020 – 2022 yang diperoleh dari data jumlah yang terdeteksi positif dan kasus aktif yang berasal dari website Corona Jakarta. Kasus positif merupakan pasien yang terkonfirmasi covid-19 dan kasus aktif merupakan pasien yang masih dalam perawatan yang dilakukan pada fasilitas kesehatan maupun secara mandiri. Pada penelitian ini menggunakan dua buah variabel data yang meliputi independen dan dependen.

1. Variabel Independen

Variabel independen merupakan variabel data-data yang digunakan untuk membuat klasterisasi yang direncanakan. Variabel ini akan memberikan pengaruh terhadap perubahan pada variabel dependen. Variabelindependent dalam penelitian ini yang terdiri dari: Positif, Dirawat, Sembuh, Meninggal dan isolasi dirumah atau isoman.

2. Variabel Dependen

Variabel dependen merupakan luaran klasterisasi yang dihasilkan dari model data mining yang dibuat, variabel ini dipengaruhi oleh variabel independent. Variabel dependen yang digunakan dalam penelitian ini yang terdiri dari: klaster merah untuk mengindikasi daerah yang rawan dengan jumlah pasien Covid-19 terbanyak. Klaster kuning untuk mengindikasi daerah yang cukup rawan dengan jumlah pasien Covid-19 hampir banyak. Klaster hijau untuk mengindikasi daerah yang aman dengan jumlah pasien Covid-19 tidak terlalu banyak.

3.6. Kerangka Kerja

Metode CRISP-DM (*Cross-Industry Standard Process for Data Mining*) merupakan model proses standar terbuka yang menggambarkan pendekatan umum yang digunakan untuk penambangan data. Tujuan dari CRISP-DM untuk menyediakan proses penambangan data yang andal dan berulang. Hal tersebut mempunyai kerangka kerja yang seragam dengan pedoman dan fleksibel untuk beradaptasi dengan masalah bisnis dan kumpulan data yang berbeda [36].

Dalam penelitian ini akan melakukan klusterisasi wilayah penyebaran Covid-19 di DKI Jakarta dengan penerapan data mining menggunakan metode *K-Means* dan *Fuzzy C-Means*, oleh karena itu untuk menerapkan teknik data mining dalam penelitian ini terdapat metode penelitian yang digunakan yaitu metode CRISP-DM. Berikut terdapat tahapan dalam CRISP-DM (*Cross-Industry Standard Process for Data Mining*) yaitu:

1. *Business Understanding*

Pada penelitian ini tujuan bisnisnya adalah memetakan klusterisasi wilayah penyebaran Covid-19 di DKI Jakarta yang diperlukan sebagai referensi menetapkan prioritas dan kebijakan penanganan pada penyebaran Covid-19 agar pemerintah setempat dapat menentukan kebijakan yang tepat dan sesuai dengan kasus di setiap daerah provinsi.

2. *Data understanding*

Pada tahap data understanding merupakan tahap persiapan yang dilakukan untuk melakukan penetapan dan validasi data yang dibutuhkan untuk pembentukan klusterisasi wilayah penyebaran Covid-19 di DKI Jakarta.

3. *Data Preparation*

Pada tahap data preparation merupakan tahap mempersiapkan data dengan menyesuaikan dataset sesuai dengan kebutuhan pada tahap pemodelan. Pada data preparation dilakukan seleksi terhadap data yang digunakan, data preparation yang bertujuan untuk mempersiapkan data mentah menjadi data yang siap untuk tahap pemodelan dan melakukan transformasi data.

4. *Modeling*

Pada tahap *modeling* merupakan tahapan pembuatan model menggunakan teknik data mining klustering, menggunakan metode *K-means Klustering* dan *fuzzy C-means Klustering* terhadap dataset penyebaran Covid-19 di DKI Jakarta. Pada tahap ini diimplementasikan dalam bentuk visualisasi menggunakan bahasa pemrograman Python dengan *tools Jupyter Notebook* dan *Microsoft Excel* untuk melakukan analisis hasil dari metode *K-means Klustering* dan *fuzzy C-means Klustering* terhadap klusterisasi penyebaran Covid-19 di DKI Jakarta.

5. *Evaluation*

Tahap *evaluation* merupakan tahap pengukuran terhadap model klusterisasi yang dihasilkan dari tahap *modeling*, dalam hal ini mengukur performansi dan akurasi dari masing-masing algoritma yang digunakan pada model tersebut yaitu *K-Means* dan *Fuzzy C Means*.

6. *Deployment*

Pada tahap *deployment* merupakan tahap implementasi hasil klusterisasi dari metode *K-Means* dan *fuzzy C-Means* pada kasus Covid-19 di DKI Jakarta yang digunakan untuk website corona Jakarta.

3.7. Teknik Analisis Data

Analisis data merupakan upaya mencari dan menata secara sistematis catatan hasil observasi, wawancara, dan lainnya untuk meningkatkan pemahaman peneliti tentang kasus yang diteliti dan menyajikannya sebagai temuan bagi orang lain. Teknik analisis data terdapat beberapa tahapan yaitu reduksi data (*Data Reduction*), penyajian data (*Data Display*) dan penarikan simpulan/verifikasi (*Conclusion drawing/verification*) [38].

1. Reduksi Data (*Data Reduction*)

Pada tahap reduksi data dilakukan dengan melakukan analisis yang menjelaskan, menggolongkan, menyederhanakan dan mengorganisasi data sehingga bisa ditarik kesimpulan dan diverifikasi. Pada tahap ini berfokus pada

tujuan yang ingin dicapai. Jika terdapat data yang dianggap tidak penting dari data kasus Covid-19 di DKI Jakarta akan dilakukan reduksi data untuk menelusuri pola kasus Covid-19 di DKI Jakarta.

2. Penyajian Data (*Data Display*)

Pada tahap penyajian data merupakan kegiatan ketika sekumpulan informasi disusun sehingga memberikan kemungkinan akan adanya penarikan kesimpulan dan pengambilan tindakan. Bentuk penyajian data dalam penelitian ini berupa visualisasi hasil klasterisasi kasus Covid-19 di DKI Jakarta yang menggunakan alat bantu bahasa pemrograman *Python*. *Python* merupakan bahasa pemrograman simpel yang digunakan untuk komputasi statistik, data mining dan grafis [39]. Dalam penelitian ini editor yang digunakan untuk menulis kode program yaitu *Jupyter Notebook*. Karena *Jupyter notebook* merupakan sebuah aplikasi untuk menulis kode bahasa pemrograman Python dalam bentuk website yang berada pada localhost komputer, setiap perintah dapat dijelaskan dalam satu halaman dan terdapat fitur yang menampilkan hasil visualisasi berupa grafik yang berada pada satu *cell* yang sama [39].

3. Penarikan simpulan/verifikasi (*Conclusion/verivication*)

Pada tahap ini yaitu penarikan kesimpulan yang didukung dengan hasil penelitian berupa klasterisasi kasus Covid-19 di DKI Jakarta dalam bentuk visualisasi.

Dari Tabel 3.3 terdapat definisi, kelebihan dan kekurangan dari *tools* yang digunakan dalam penelitian ini. Dengan menggunakan *Jupyter Notebook* dapat memvisualisasi data dalam bentuk grafik dan dapat menjelaskan setiap perintah dari setiap bagian program pemodelan yang dibuat [39].

Tabel 3.3 Perbandingan *Tools Data Mining* Yang Digunakan

<i>Tools</i>	Definisi	Kelebihan	Kekurangan
<i>Jupyter Notebook</i>	Alat <i>open source</i>	• Visualisasi data,	• <i>Jupyter Notebook</i>

	<p>berbasis Python yang berinteraksi langsung dengan data para ilmuwan dikarenakan reatif mudah dipelajarinya.</p>	<p>dapat melakukan visualiasasi data dalam bentuk grafik, tableau.</p> <ul style="list-style-type: none"> • Berbagi kode, dapat melakukan berbagi kode secara <i>cloud</i> atau melalui jaringan internet untuk melihat kode, menjalankan kode dan menampilkannya pada browser, seperti GitHub • Interaksi langsung dengan kode, kode dalam <i>jupyter</i> tidak statis sehingga dapat diedit dan dijalankan kembali dan menghasilkan umpan balik pada browser. • Mendokumentasikan kode, dalam <i>jupyter</i> dapat menjelaskan secara bertahap pada potongan kode yang ingin dijelaskan. 	<p>tidak mandiri yaitu membutuhkan <i>runtime jupyter</i>.</p> <ul style="list-style-type: none"> • Status sesi tidak dapat disimpan dengan mudah yaitu kode yang dijalankan tidak dapat dipertahankan dan dikembalikan sehingga setiap membuat notebook harus dijalankan kembali kodenya untuk memulihkan kondisinya. • Tidak ada <i>debuging</i> interaktif atau IDE lain
<p>Google Colaboratory</p>	<p>Perangkat lunak yang dapat digunakan untuk mengeksekusi kode python melalui <i>web browser</i> [40].</p>	<ul style="list-style-type: none"> • Dapat digunakan untuk <i>machine learning</i>, analisis data dan pendidikan. • Layanan mirip <i>jupyter notebook</i> yang telah dihosting dan disediakan oleh 	<ul style="list-style-type: none"> • <i>Resource</i> yang disediakan tidak terjamin dan terbatas yang memiliki batas penggunaan <i>resource</i> tersebut [40]. • Tidak dapat menjamin jenis

		google yang memiliki akses bebas dan <i>resource</i> gratis [40].	<i>hardware</i> yang diperoleh atau kapasitas memori yang diperlukan [41].
--	--	---	--

Dari tabel 3.3 merupakan perbandingan *tools* pada data mining, dimana setiap *tools* memiliki kelebihan dan kekurangannya. Dalam penelitian yang dilakukan, peneliti menggunakan *tools Jupyter Notebook* karena pada *Jupyter* terdapat berbagai teknik dalam data mining yaitu teknik *supervised* maupun *unsupervised learning*. Dalam *tools* tersebut juga menyediakan *library* yang dapat digunakan dalam pembuatan sistem, *library* merupakan kumpulan modul yang berisi kumpulan kode yang dapat digunakan berulang kali dalam program yang berbeda [18]. Dalam *library* pada Python yang tersedia diantaranya yaitu numPy, Pandas, matplotlib, sklearn, serta menyediakan berbagai algoritma dan menghasilkan *output* berupa visualisasi dalam bentuk grafis [18]. Selain itu, *tools* ini sangat berguna di bidang *data science* dan *machine learning* mengingat penelitian yang dilakukan menggunakan metode klusterisasi yang akan diimplementasikan kedalam bahasa pemrograman Python, karena *tools* ini terintegrasi dengan Python dan *library* yang ada tidak mempunyai batasan waktu. Hal ini berbeda dengan *google laboratory* yang memiliki *resource* yang tidak terjamin dan terbatas serta tidak dapat menjamin kapasitas memori yang diperlukan.