

**Pembelajaran Mesin**  
**Case Based 2 - Unsupervised Learning**  
**Clustering**



Disusun oleh:

Rayhan Suryatama Raharyawhedi

1301204435

**Universitas Telkom**

**Jl. Telekomunikasi. 1, Terusan Buahbatu - Bojongsoang, Telkom University,  
Sukapura, Kec. Dayeuhkolot, Kabupaten Bandung, Jawa Barat 40257**

Tugas ini dikerjakan secara mandiri. Semua referensi terkait pembangunan tugas ini terlampir pada bab Referensi. Penulis sudah berusaha semaksimal mungkin untuk membuat laporan ini orisinal. Segala bentuk plagiarisme pada laporan ini mungkin dikarenakan kesamaan dari referensi tercantum atau kesalahan yang tidak disengaja. Tugas ini juga dikerjakan dengan tidak melanggar ketentuan perkuliahan dan kode etik akademisi.

## Daftar Isi

1. Rumusan Masalah .....	4
2. Analisis Data.....	4
2.1. Dimensi Dataset.....	4
2.2. Kolom Tidak Berguna/Tidak Dapat Diproses.....	5
2.3. Nilai Kosong .....	5
2.4. Rentang Nilai .....	6
3. Pra-pemrosesan Data.....	6
3.1. Menghilangkan Nilai Kosong.....	7
3.2. Scaling Data.....	7
3.3. Reduksi Dimensi.....	8
4. Penerapan Algoritma Clustering.....	9
5. Evaluasi.....	13
6. Presentasi Video .....	14
7. Link Google Colab .....	14
8. Referensi .....	14

## 1. Rumusan Masalah

Permasalahan yang disinggung pada laporan tugas Case Based 2 ini adalah pengelompokan/pelabelan data kepada sebuah dataset water-treatment.data (<https://archive.ics.uci.edu/ml/datasets/Water+Treatment+Plant>). Dataset tersebut yang belum memiliki label akan diolah sedemikian rupa agar memiliki label yang nantinya dapat diolah lebih lanjut. Proses pelabelan data untuk Case Based 2 ini akan menggunakan metode K-Means, dimana nanti data akan dikelompokkan berdasarkan banyak kelompok K yang telah ditentukan. Data dikelompokkan berdasarkan nilai tengah yang nantinya akan dibuat dan dikalkulasi.

## 2. Analisis Data

Dataset water-treatment akan dianalisis kualitasnya yang nantinya dari kualitas tersebut akan dibenahi pada saat pra-pemrosesan data. Pengecekan kualitas tersebut meliputi:

- Dimensi dataset
- Kolom yang tidak berguna/tidak dapat diproses
- Nilai kosong
- Rentang nilai

### 2.1. Dimensi Dataset

Dataset water-treatment memiliki dimensi yang cukup besar. Dengan adanya 39 kolom, maka dapat dibilang dataset ini berdimensi 39.

	ID	Q-E	ZN-E	PH-E	DBO-E	DQO-E	SS-E	SSV-E	SED-E	COND-E	...	COND-S	RD-DBO-P	RD-SS-P	RD-SED-P	RD-DBO-S	RD-DQO-S	RD-DBO-G	RD-DQO-G	RD-SS-G	RD-SED-G
0	D-1/3/90	44101.0	1.50	7.8	?	407	166.0	66.3	4.5	2110	...	2000.0	?	58.8	95.5	?	70.0	?	79.4	87.3	99.6
1	D-2/3/90	39024.0	3.00	7.7	?	443	214.0	69.2	6.5	2660	...	2590.0	?	60.7	94.8	?	80.8	?	79.5	92.1	100
2	D-4/3/90	32229.0	5.00	7.6	?	528	186.0	69.9	3.4	1666	...	1888.0	?	58.2	95.6	?	52.9	?	75.8	88.7	98.5
3	D-5/3/90	35023.0	3.50	7.9	205	588	192.0	65.6	4.5	2430	...	1840.0	33.1	64.2	95.3	87.3	72.3	90.2	82.3	89.6	100
4	D-6/3/90	36924.0	1.50	8.0	242	496	176.0	64.8	4.0	2110	...	2120.0	?	62.7	95.6	?	71.0	92.1	78.2	87.5	99.5
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
522	D-26/8/91	32723.0	0.16	7.7	93	252	176.0	56.8	2.3	894	...	942.0	?	62.3	93.3	69.8	75.9	79.6	78.6	96.6	99.6
523	D-27/8/91	33536.0	0.32	7.8	192	346	172.0	68.6	4.0	988	...	950.0	?	58.3	97.8	83.0	59.1	91.1	74.6	90.7	100
524	D-28/8/91	32922.0	0.30	7.4	139	367	180.0	64.4	3.0	1060	...	1136.0	?	65.0	97.1	76.2	66.4	82.0	77.1	88.9	99
525	D-29/8/91	32190.0	0.30	7.3	200	545	258.0	65.1	4.0	1260	...	1326.0	39.8	65.9	97.1	81.7	70.9	89.5	87.0	89.5	99.8
526	D-30/8/91	30488.0	0.21	7.5	152	300	132.0	69.7	?	1073	...	1224.0	?	69.5	?	81.7	76.4	?	81.7	86.4	?

527 rows x 39 columns

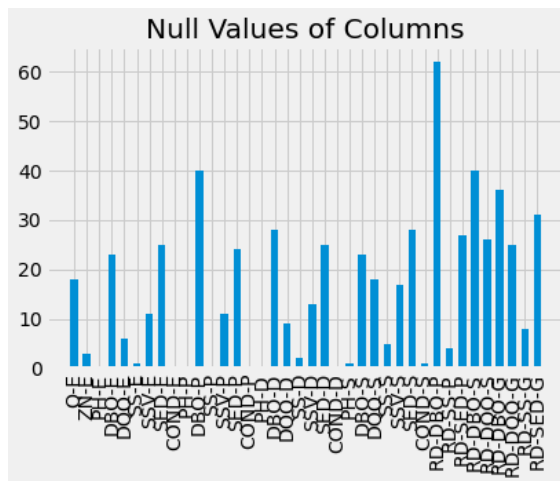
527 rows x 39 columns

## 2.2.Kolom Tidak Berguna/Tidak Dapat Diproses

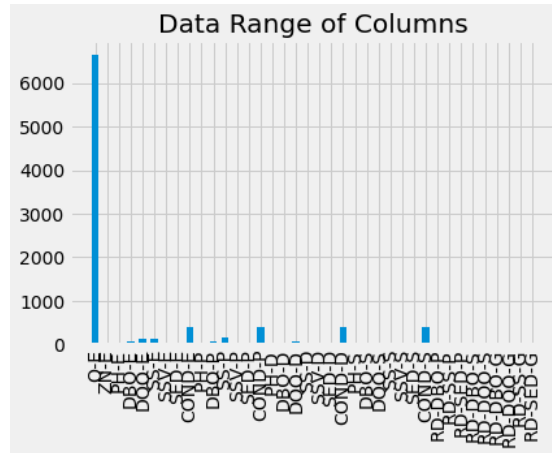
ID	
0	D-1/3/90
1	D-2/3/90
2	D-4/3/90
3	D-5/3/90
4	D-6/3/90
...	...
522	D-26/8/91
523	D-27/8/91
524	D-28/8/91
525	D-29/8/91
526	D-30/8/91

Dari analisis lebih lanjut, kolom ‘ID’ tidak dibutuhkan karena kolom tersebut, pertama, bernilai string/object, dan kedua, tidak berpengaruh signifikan ke data-data lain. ‘ID’ hanya dibutuhkan untuk pemberian identifikasi unik untuk setiap datanya. Index dari baris sudah cukup merepresentasikan hal tersebut. Maka kolom ‘ID’ tidak akan digunakan untuk proses Clustering ini.

## 2.3.Nilai Kosong



## 2.4. Rentang Nilai



Rentang nilai ini diambil dari perbandingan standar deviasi dari setiap kolom. Rentang nilai dari dataset water-treatment cukup kecil, namun ada satu kolom yang memiliki rentang nilai tertinggi, yaitu 'O-E' dikarenakan kolom tersebut adalah kolom utama pada dataset ini, yang berisikan data-data penting dari kolom lain. Perbandingan nilai yang jauh ini akan diselesaikan masalahnya dengan menggunakan metode scaling Min Max Scaling.

## 3. Pra-pemrosesan Data

Langkah setelah melakukan analisis data, hasil analisis tersebut akan dibenahi pada tahap Pra-Pemrosesan data ini. Data-data yang sekiranya tidak dibutuhkan atau masih kurang tepat jika digunakan akan diselesaikan masalahnya agar proses clustering nanti dapat berjalan dengan baik. Pra-pemrosesan data meliputi:

- Menghilangkan nilai kosong
- Scaling data
- Reduksi dimensi

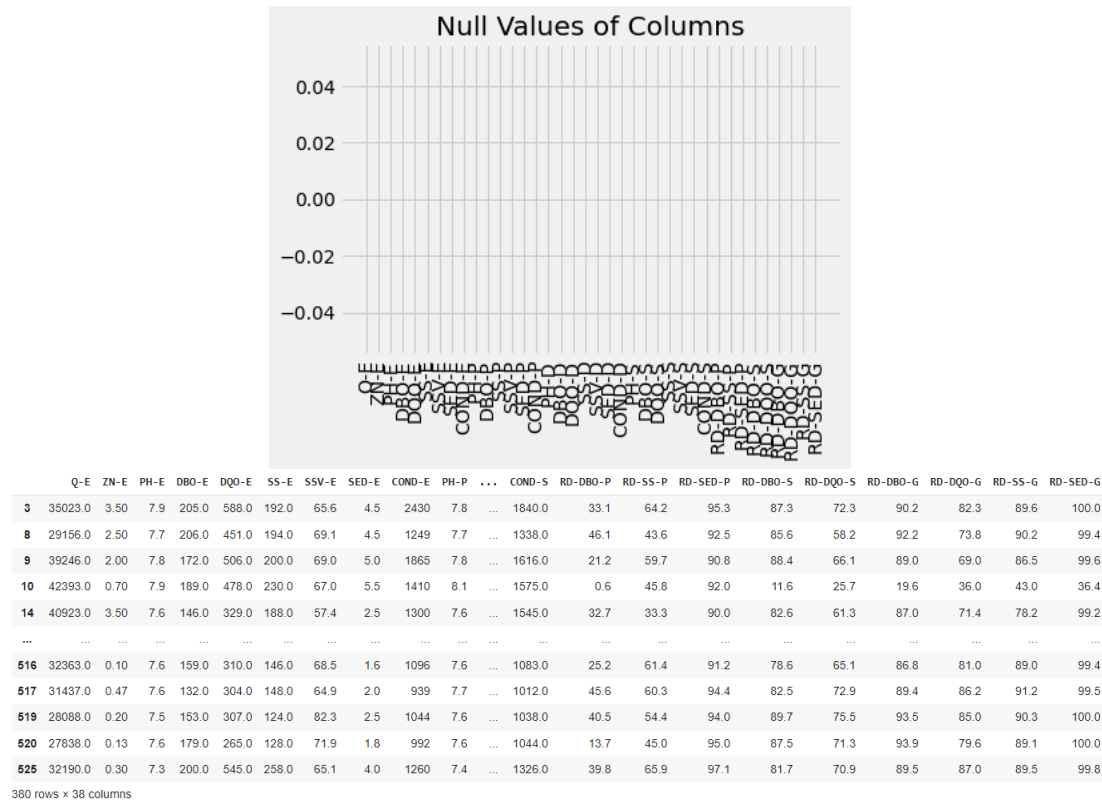
### 3.1.Menghilangkan Nilai Kosong

Nilai kosong yang telah dianalisis di bab sebelumnya akan ditangani dengan cara menghapus baris yang memiliki nilai kosong pada kolom tertentu.

Snippet Kode:

```
df_main = df_main.dropna(axis=0)
```

Hasil:



380 rows x 38 columns

Setelah dilakukan penghapusan baris yang memiliki nilai kosong, dapat dilihat sudah tidak ada nilai kosong di semua kolom, walaupun total jumlah baris dataset berkurang.

### 3.2.Scaling Data

Scaling dataset ini akan menggunakan cara Min Max Scaler yang memiliki rumus:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

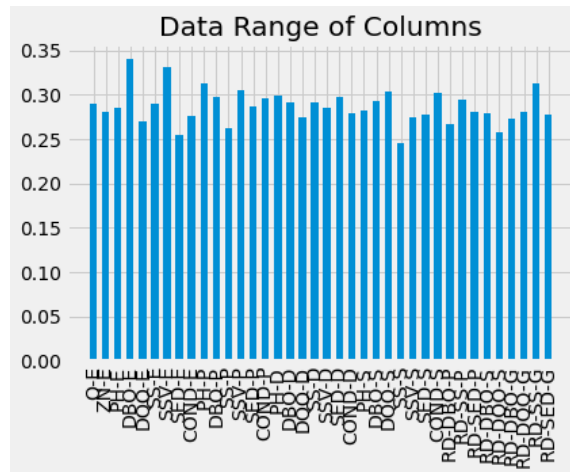
Pada proses scaling juga digunakan sebuah library dari scikit-learn yaitu dari kelas preprocessing-nya, MinMaxScaler().

Snippet Kode:

```
def dataScaling(df):  
    scaler = MinMaxScaler()  
  
    arr_scaled = scaler.fit_transform(df_main)  
  
    return arr_scaled
```

```
data_scaled = dataScaling(df_main)
```

Hasil:



Hasil scaling data memperlihatkan bahwa data-data yang telah dilakukan proses scaling ini berkurang rentang data-nya. Sekarang data hanya berada pada rentang 0 hingga 1.

### 3.3.Reduksi Dimensi

Reduksi dimensi dilakukan agar data yang akan diproses nanti dapat divisualisasikan dengan benar. Karena dimensi data dari water-treatment yang memiliki 38 (39 sebelum dilakukan pembuangan kolom tidak berguna) maka perlu dilakukan reduksi dimensi menjadi 2 dimensi saja agar visualisasi dari data nanti akan lebih jelas.

Reduksi dimensi ini mengambil 2 kolom utama yang diambil dari kumpulan kolom-kolom lain.

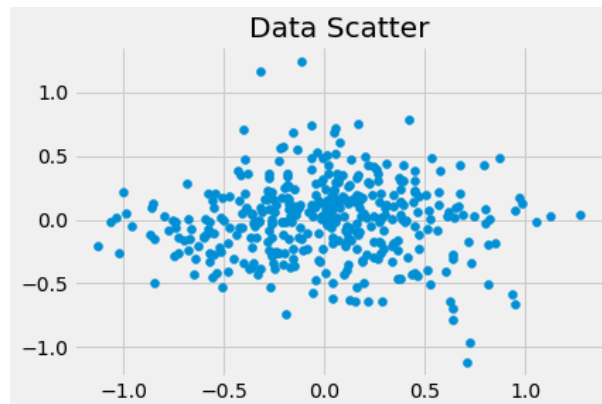
Proses reduksi dimensi menggunakan library scikit-learn, dari kelas decomposition, PCA().

Snippet Kode:

```
pca = PCA(n_components=2)

pc = pca.fit_transform(data_scaled)
```

Hasil:





#### 4. Penerapan Algoritma Clustering

Algoritma Clustering yang akan digunakan pada tugas Case-Based 2 ini adalah K-Means. K-Means ini bekerja dengan cara menentukan berapa K cluster yang diinginkan, kemudian membuat sebuah centroid (data tengah) berjumlah K baru yang diambil dari titik random dari data, setelah itu mengelompokkan semua data dengan cluster centroid terdekat. Setelah pengelompokkan pertama selesai akan dilakukan penentuan centroid baru yang kemudian data-data akan dikelompokkan ulang, hingga nilai dari cluster tidak berubah lagi.

Pseudocode K-Means Clustering:

```
KMeansClustering
Algoritma
    init K centroids
    repeat
        pairing closest point to K cluster
        recalculate centroids
    until Centroids don't change
```

Snippet Kode Model K-Means:

\*terlampir di halaman selanjutnya

```

1 class KMeansClustering:
2     def __init__(self, k, max_iter=100, random_state=42):
3         self.k = k
4         self.max_iter = max_iter
5         self.random_state = random_state
6
7     def initCentroids(self, data):
8         np.random.RandomState(self.random_state)
9         rand_i = np.random.permutation(data.shape[0])
10        centroids = data[rand_i[:self.k]]
11
12        return centroids
13
14    def computeCentroids(self, data, labels):
15        centroids = np.zeros((self.k, data.shape[1]))
16        for i in range(self.k):
17            centroids[i, :] = np.mean(data[labels == i, :], axis=0)
18
19        return centroids
20
21    def computeDistance(self, data, centroids):
22        dist = np.zeros((data.shape[0], self.k))
23        for i in range(self.k):
24            row_norm = np.linalg.norm(data - centroids[i, :], axis=1)
25            dist[:, i] = np.square(row_norm)
26
27        return dist
28
29    def closestToCluster(self, dist):
30        return np.argmin(dist, axis=1)
31
32    def computeSSE(self, data, labels, centroids):
33        dist = np.zeros(data.shape[0])
34        for i in range(self.k):
35            dist[labels == i] = np.linalg.norm(data[labels == i] - centroids[i], axis=1)
36
37        return np.sum(np.square(dist))
38
39    def fit(self, data):
40        self.centroids = self.initCentroids(data)
41        for i in range(self.max_iter):
42            temp_centroids = self.centroids
43            dist = self.computeDistance(data, temp_centroids)
44            self.labels = self.closestToCluster(dist)
45            self.centroids = self.computeCentroids(data, self.labels)
46            if np.all(temp_centroids == self.centroids):
47                break
48
49        self.error = self.computeSSE(data, self.labels, self.centroids)
50
51        # def predict(self, data):
52        #     dist = self.computeDistance(data, temp_centroids)
53        #     return self.closestToClust(dist)

```

Kode mereferensi sebuah artikel dari @afrizalfir di medium.com

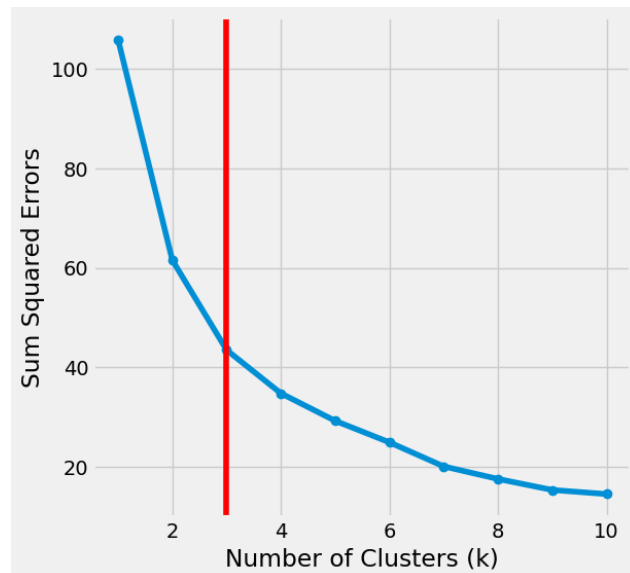
(<https://medium.com/@afrizalfir/kmeans-clustering-dan-implementasinya-5e967dc604cf>)

Setelah perancangan model K-Means selesai, dilanjutkan dengan Elbow Method untuk mencari K cluster paling optimal.

Snippet Kode Elbow Method:

```
1 sse = []
2 for i in range(1, 11):
3     kmeans = KMeansClustering(k=i, max_iter=100)
4     kmeans.fit(pc)
5     sse.append(kmeans.error)
6
7 plt.figure(figsize=(6,6))
8 plt.plot(range(1, 11), sse, '-o')
9 plt.xlabel(r'Number of Clusters (k)')
10 plt.ylabel('Sum Squared Errors')
11 plt.axvline(x=3, color='r')
12 plt.show
```

Hasil:



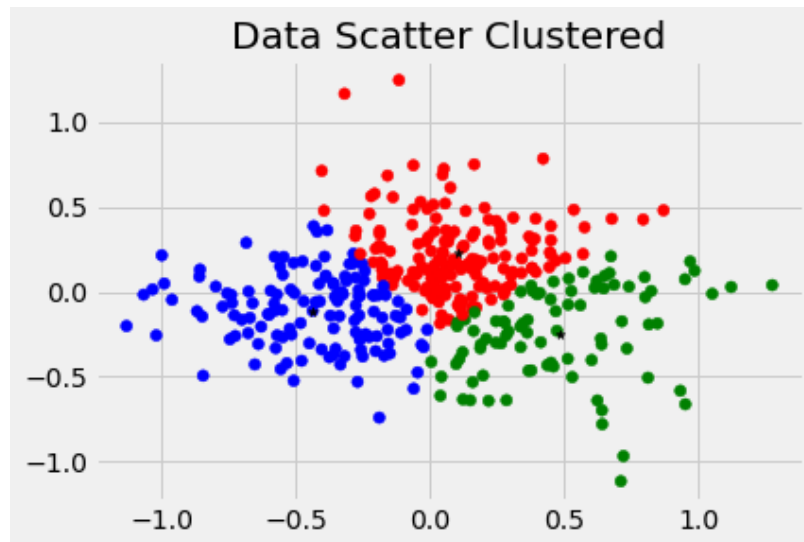
Dari analisis Elbow Method yang melihat Sum Square Error dari setiap K cluster, dapat dilihat bahwa K = 3 adalah jumlah cluster paling optimal karena K = 3 ini memiliki penurunan nilai Sum Square Error yang mulai landai daripada K = 4 dan seterusnya.

Dengan menggunakan K = 3 yang didapat dari Elbow Method, maka dapat dibuat sebuah cluster menggunakan K-Means dengan nilai K = 3.

Snippet Kode Implementasi K-Means:

```
1 kmeans = KMeansClustering(k=3)
2 kmeans.fit(pc)
3
4 colors = {
5     0:'r',
6     1:'b',
7     2:'g'
8 }
9 kmeans_labels = [colors[label] for label in kmeans.labels]
10
11 plt.scatter(pc[:, 0], pc[:, 1], c=kmeans_labels)
12 plt.scatter(kmeans.centroids[:, 0], kmeans.centroids[:, 1], marker='*', s=24, c='k')
13 plt.title('Data Scatter Clustered')
14 plt.show()
```

Hasil:



Setelah dilakukan proses implementasi model K-Means yang telah dibuat, dapat dilihat bahwa data-data yang dulunya tidak memiliki label kelompok sudah memiliki kelompok. Juga terlihat centroid mana yang menjadi patokan kelompok tersebut (centroid ditandai dengan simbol bintang (★) hitam pada graf).

## **5. Evaluasi**

Proses K-Means ini mudah digunakan, namun kurang cocok untuk beberapa dataset yang tidak sesuai kriteria K-Means. Kriteria yang dimaksud seperti dimensi data, dimensi paling cocok untuk pengelompokkan K-Means agar visualisasi data-nya mudah dipahami adalah data 2 dimensi saja. Kriteria lain adalah bentuk data. Karena bentuk data dari water-treatment ini berkumpul di satu titik, K-Means cenderung cocok diterapkan, namun jika persebaran data-nya berbeda mungkin K-Means tidak akan cocok untuk data tersebut.

Cluster K dari K-Means juga tidak semena-mena ditentukan. Digunakan elbow method untuk menentukan nilai K tersebut memiliki alasan agar cluster yang dibangun adalah cluster paling optimal. Walaupun K-Means juga bisa membuat K tersendiri sesuai yang diinginkan pengolah data.

## **6. Presentasi Video**

[https://drive.google.com/file/d/14Q5xDFpa1wjX0fTm0lvAXwnm4PJVJ2NB/view?usp=share\\_link](https://drive.google.com/file/d/14Q5xDFpa1wjX0fTm0lvAXwnm4PJVJ2NB/view?usp=share_link)

## **7. Link Google Colab**

<https://colab.research.google.com/drive/10rU14PS-7JUWTbK0AWmtBDxtjvVenswu?usp=sharing>

## **8. Referensi**

<https://www.bradleysawler.com/engineering/ml-clustering-of-a-waste-water-treatment-plant/>

<https://medium.com/@afrizalfir/kmeans-clustering-dan-implementasinya-5e967dc604cf>

<https://predictivehacks.com/k-means-elbow-method-code-for-python/>

<https://realpython.com/k-means-clustering-python/>

<https://lms.telkomuniversity.ac.id/mod/page/view.php?id=2589254>