

Clustering

Hendrik Santoso Sugiarto

IBDA2032 – *Artificial Intelligence*

Capaian Pembelajaran

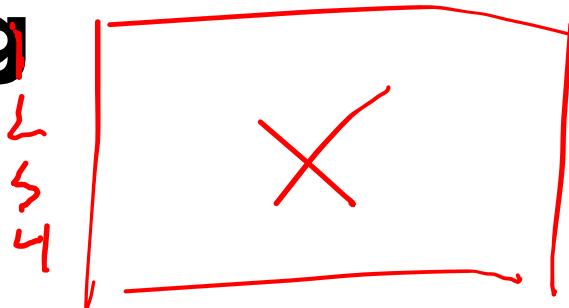
- Konsep Clustering
- Hard Clustering
- Soft Clustering

Clustering

Unsupervised Learning

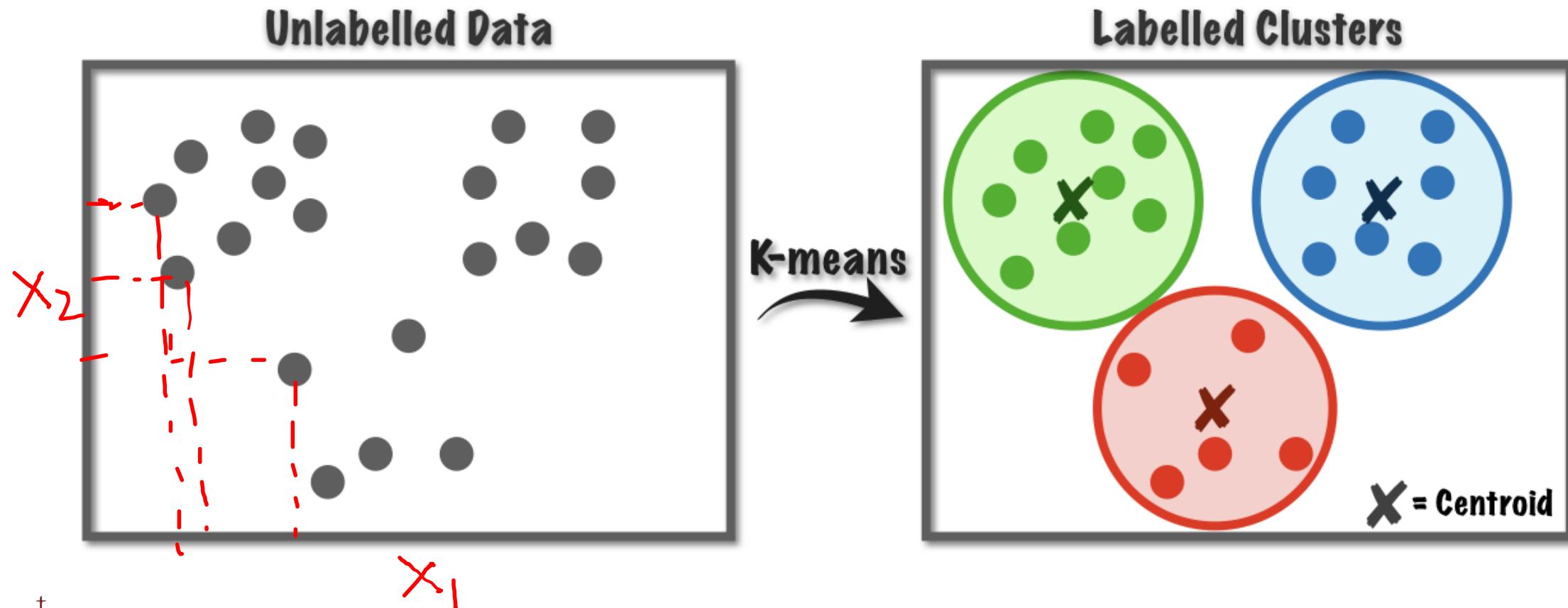
- Bentuk formal:

- Input: $x \in \mathcal{X} \in \mathbb{R}^n$
- Output: $y \in \mathcal{Y} \in \begin{cases} \{1, 2, \dots, K\} \rightarrow \text{clustering } K \text{ grup} \\ \mathbb{R}^K \rightarrow \text{embedding } K \text{ dimension} \end{cases}$
- Target function: $f: \mathcal{X} \rightarrow \mathcal{Y}$ (*unknown*)
- Training data: $D = \{(x^{(1)}), (x^{(2)}), \dots, (x^{(m)})\}$
- Hypothesis: $h: \mathcal{X} \rightarrow \mathcal{Y}$
- Hypothesis space: $h \in \mathcal{H}$



Clustering

- Menemukan struktur tersembunyi pada data
- Unsupervised learning: tidak ada informasi mengenai label



Contoh Aplikasi Clustering: Search Result

company | products | solutions | customers | demos | partners | press

jaguar

Search

Clustering by Vivisimo

Other demos | Help! | Tell us what you think!

Clustered Results

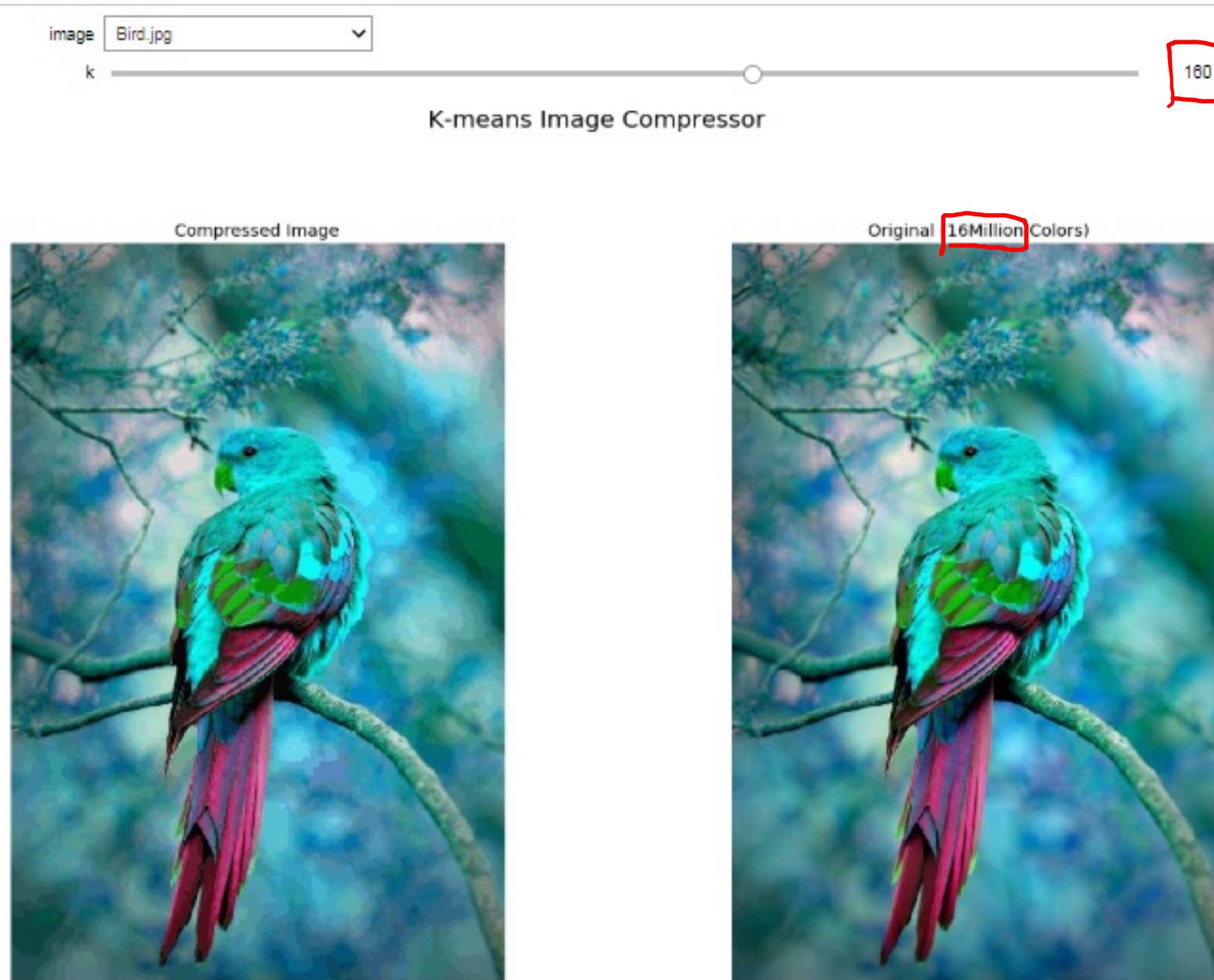
Top 185 results retrieved for the query **jaguar** (Details)

- Jaguar Cars** [new window] [frame] [preview]
Official worldwide web site of **Jaguar** Cars. Gama actual, concesionarios, historia, noticias, anuncios y servicios fina
URL: www.jaguar.com - show in clusters
Sources: Lycos 1
- Jaguar Cars** [new window] [frame] [preview]
URL: www.jaguarcars.com - show in clusters
Sources: Lycos 2, Lycos 59, Lycos 90, Lycos 97, Lycos 99
- www.jaguar-racing.com [new window] [frame] [preview]
URL: www.jaguar-racing.com - show in clusters
Sources: Lycos 3, Lycos 93, Lycos 115
- Jaguar Cars** [new window] [frame] [preview]
United States United Kingdom Germany Japan France Italy Spain...
URL: www.jaguarvehicles.com - show in clusters
Sources: Lycos 4, Lycos 8, Lycos 41, Lycos 102, Lycos 188
- Apple - Mac OS X** [new window] [frame] [preview]
... queries to find your stuff, refining the list as you narrow options. Sure you could quantify that as up to six times faster than **Jaguar**, but you'll probably think Panthers done almost before you...
URL: www.apple.com/macosx - show in clusters
Sources: Lycos 5

Find in clusters:

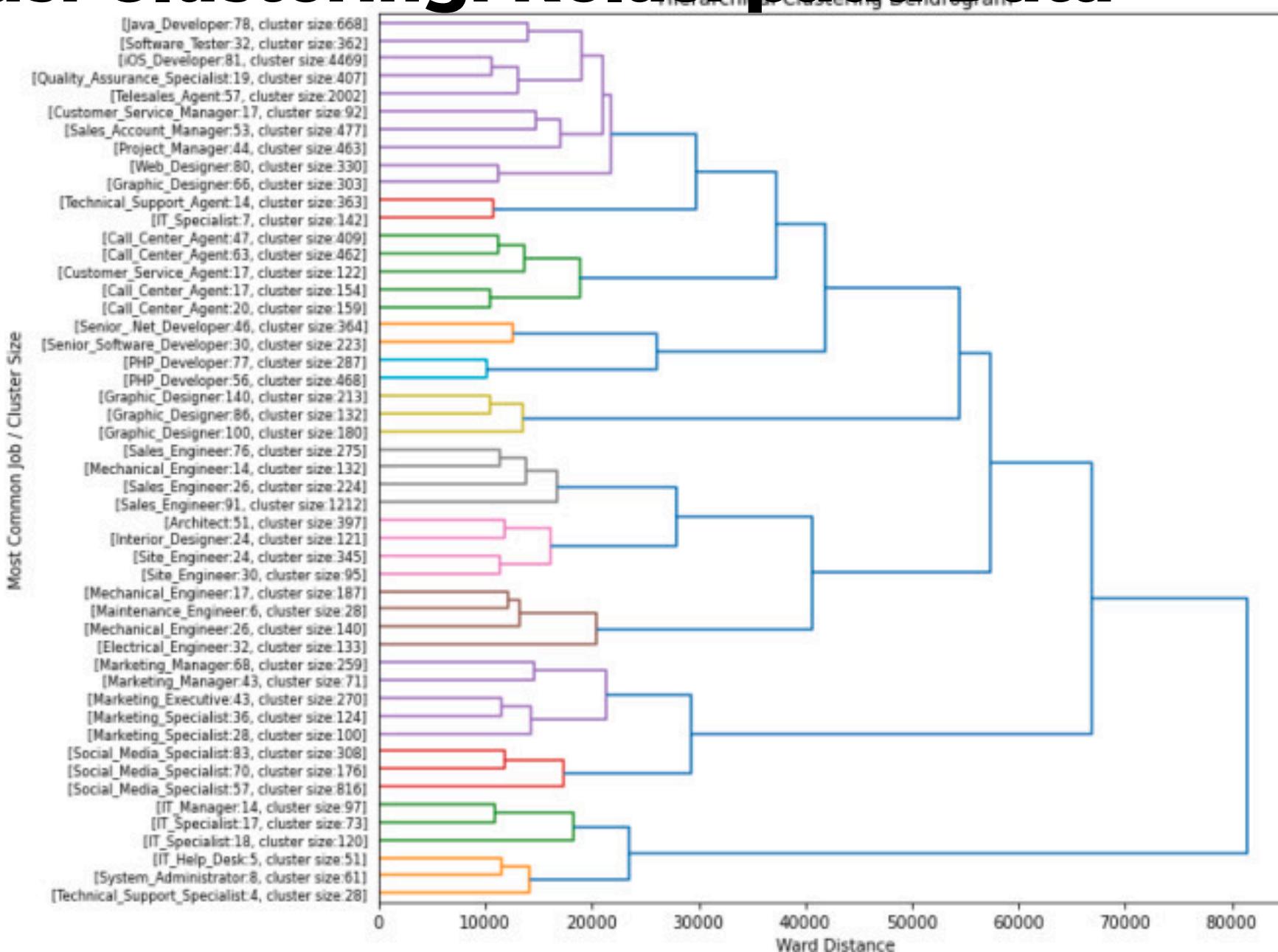
Enter Keywords

Contoh aplikasi clustering: Image Compression



Contoh Aplikasi Clustering: Relasi pada Data

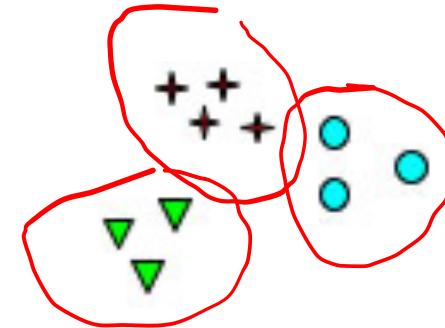
- www.sciencedirect.com/science/article/abs/pii/S1877750322002149



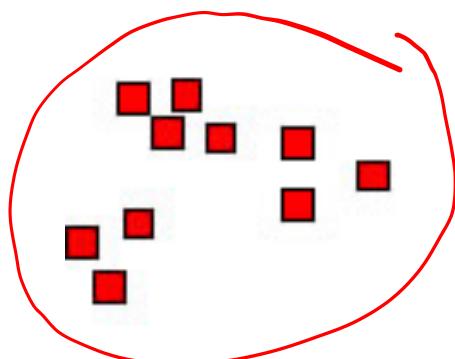
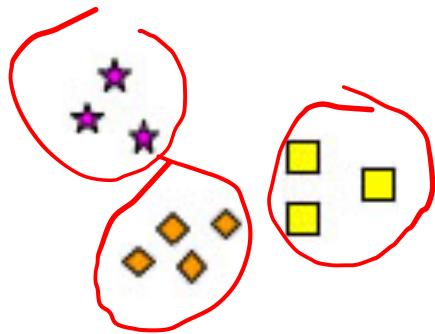
Cluster bersifat ambigu



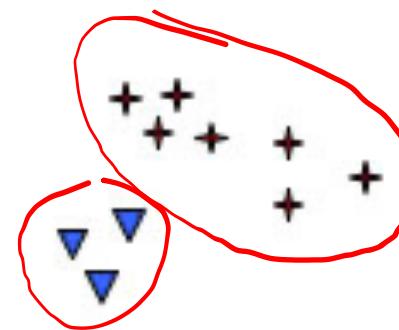
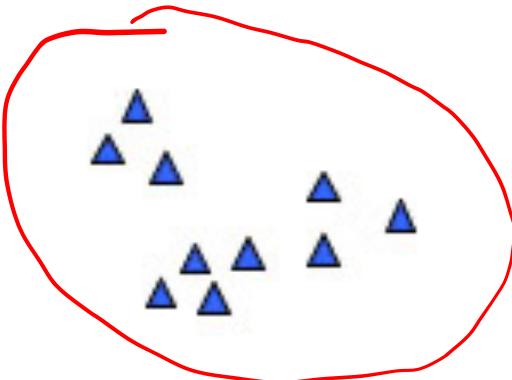
How many clusters?



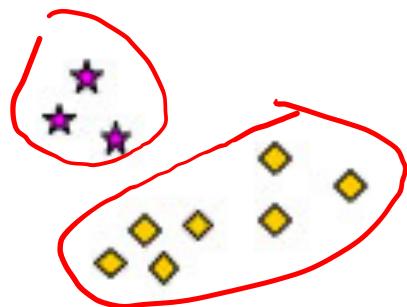
Six Clusters



Two Clusters



Four Clusters



Jenis Clustering

- Hard Clustering:
 - Partitional clustering: pembagian data menjadi subset yang tidak overlap sehingga setiap data terletak pada suatu cluster
 - Hierarchical clustering: himpunan nested cluster yang disusun menurut sebuah hierarchical tree
 - Spectral clustering: pembagian data menjadi subset yang mengikuti struktur lokal
- Soft Clustering:
 - Mixture clustering: pembagian data menjadi subset yang saling overlap sehingga setiap data dapat terletak pada beberapa cluster sekaligus

Hard Clustering

K-means Clustering

- Konsep: minimisasi total jarak intra-cluster

- Cost function:

$K = \text{Jumlah cluster}$
 $M = \text{Jumlah data}$

$$J(\boldsymbol{\theta}) = \sum_{i=1}^M \sum_{k=1}^K A_{ik} \left\| \boldsymbol{\theta}_k - \mathbf{x}^{(i)} \right\|^2$$

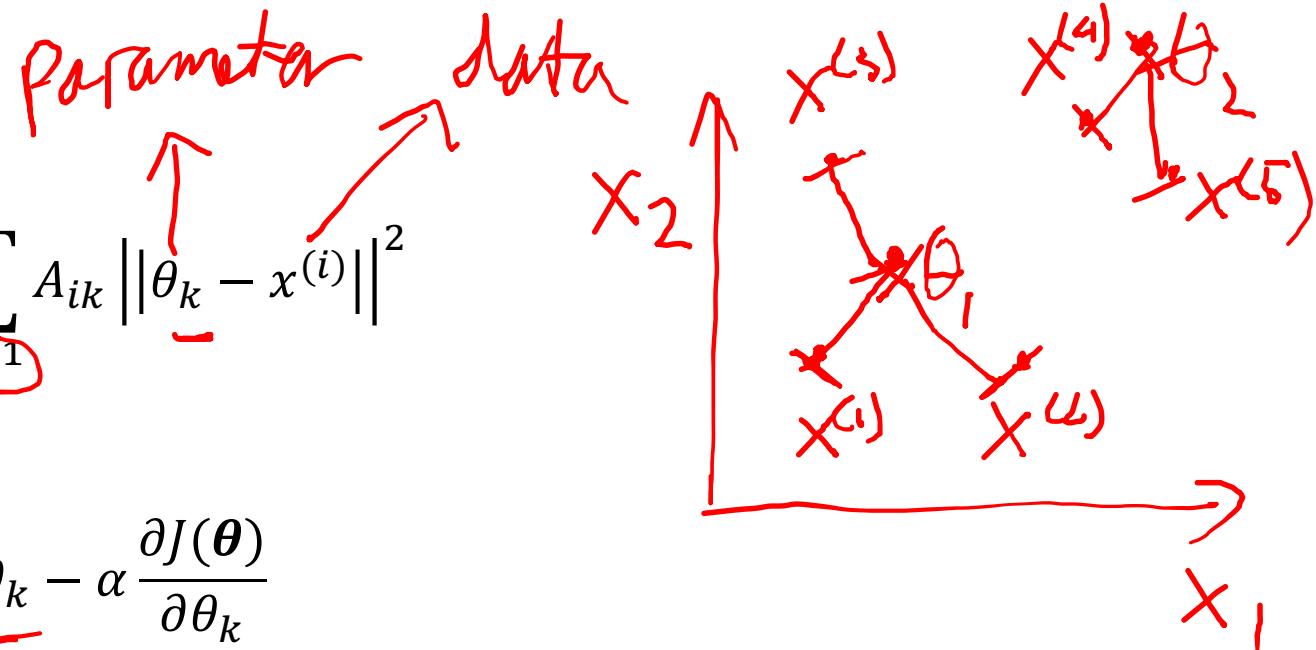
- Dimana $A_{ik} \in \{0,1\} \forall i, k$ dan $\sum_{k=1}^K A_{ik} = 1 \forall i$

- Gradient descent:

$$\boldsymbol{\theta}_k := \boldsymbol{\theta}_k - \alpha \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_k}$$

- Solusi analitik:

$$0 = \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_k} = 2 \sum_{i=1}^M A_{ik} (\boldsymbol{\theta}_k - \mathbf{x}^{(i)}) \rightarrow \sum_{i=1}^M A_{ik} (\boldsymbol{\theta}_k) = \sum_{i=1}^M A_{ik} (\mathbf{x}^{(i)})$$
$$\boldsymbol{\theta}_k = \frac{\sum_{i=1}^M A_{ik} (\mathbf{x}^{(i)})}{\sum_{i=1}^M A_{ik}}$$

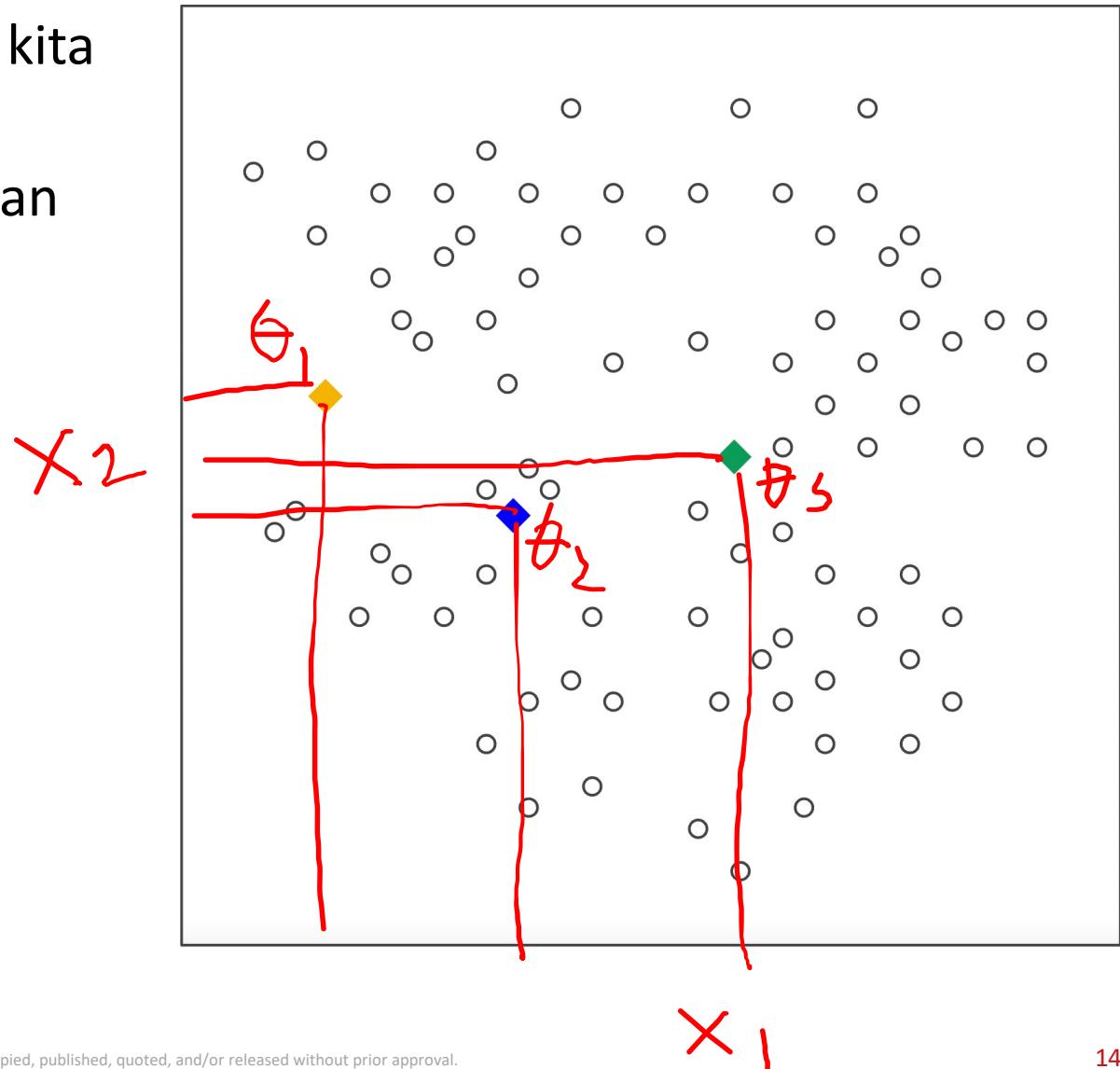


Algoritma K-Means

- <https://developers.google.com/machine-learning/clustering/algorithm/run-algorithm>
- 0. Tentukan jumlah cluster K
- 1. Algoritma memilih secara acak centroid (pusat) dari setiap cluster
- 2. Algoritma menetapkan setiap data pada centroid terdekat untuk memperoleh k cluster awal
- 3. Untuk setiap cluster, algoritma akan menghitung ulang centroid sebagai rata-rata dari setiap data cluster tersebut means
- 4. Algoritma mengulangi perhitungan centroid dan penetapan cluster sampai konvergen

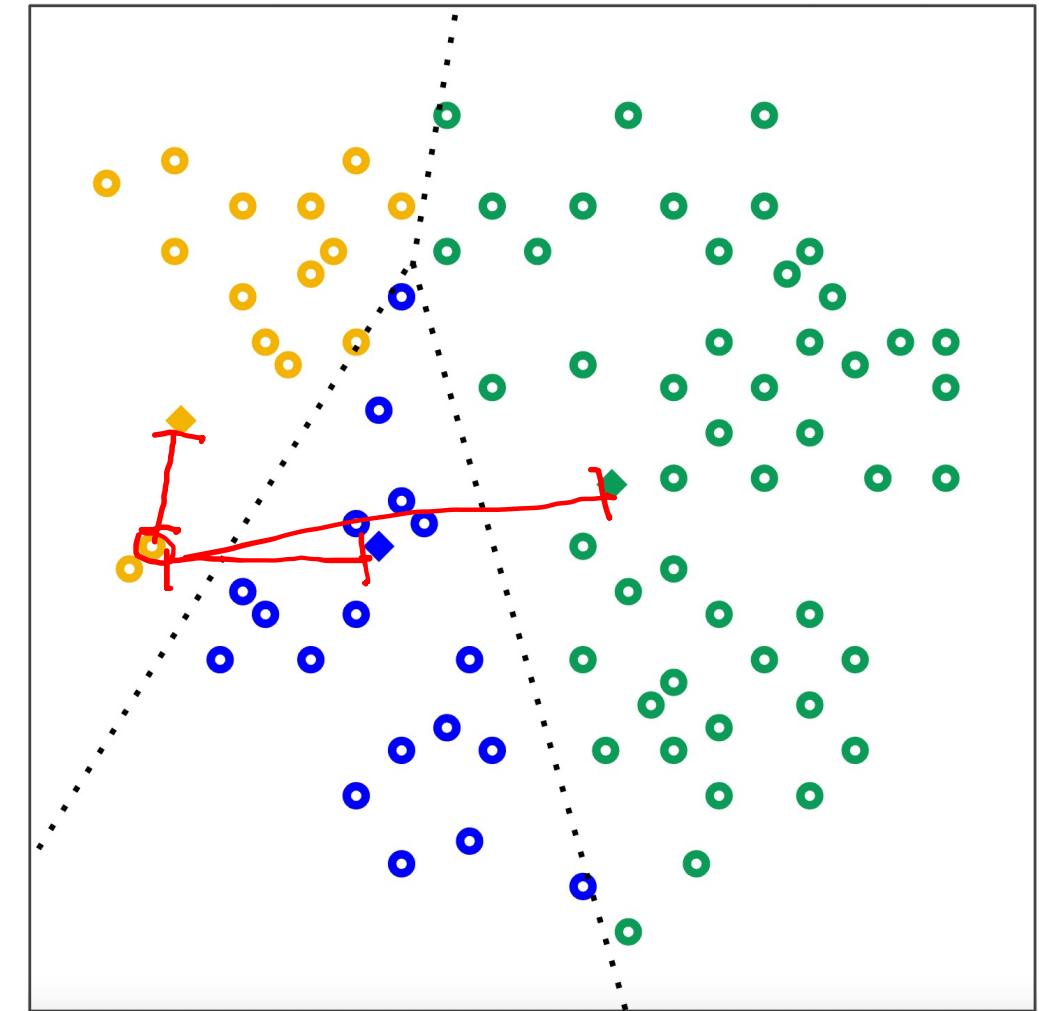
Pilih acak centroid

- Tentukan jumlah cluster (dalam kasus ini kita tentukan ada 3 clusters) $k=3$
- Lalu pilih acak posisi centroid pada sebaran data (centroid berada pada posisi yang berwarna kuning, biru, hijau)



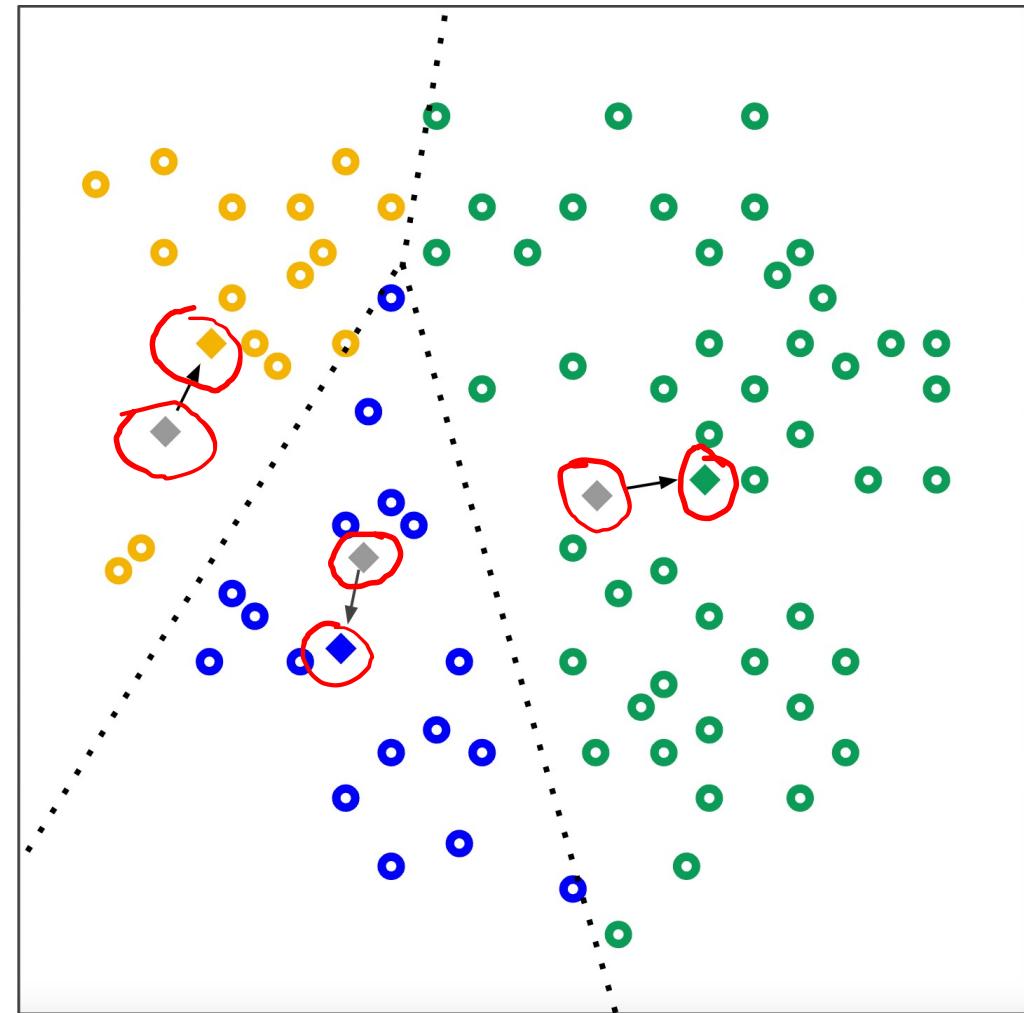
Penentuan cluster sesuai jarak centroid terdekat

- Untuk setiap data, tentukan anggota clusternya dengan cara menghitung jaraknya dengan setiap centroid
- Ialu pilih cluster dengan centroid terdekat (setiap data digolongkan menuju cluster kuning, hijau, dan biru menurut kedekatan lokasinya)



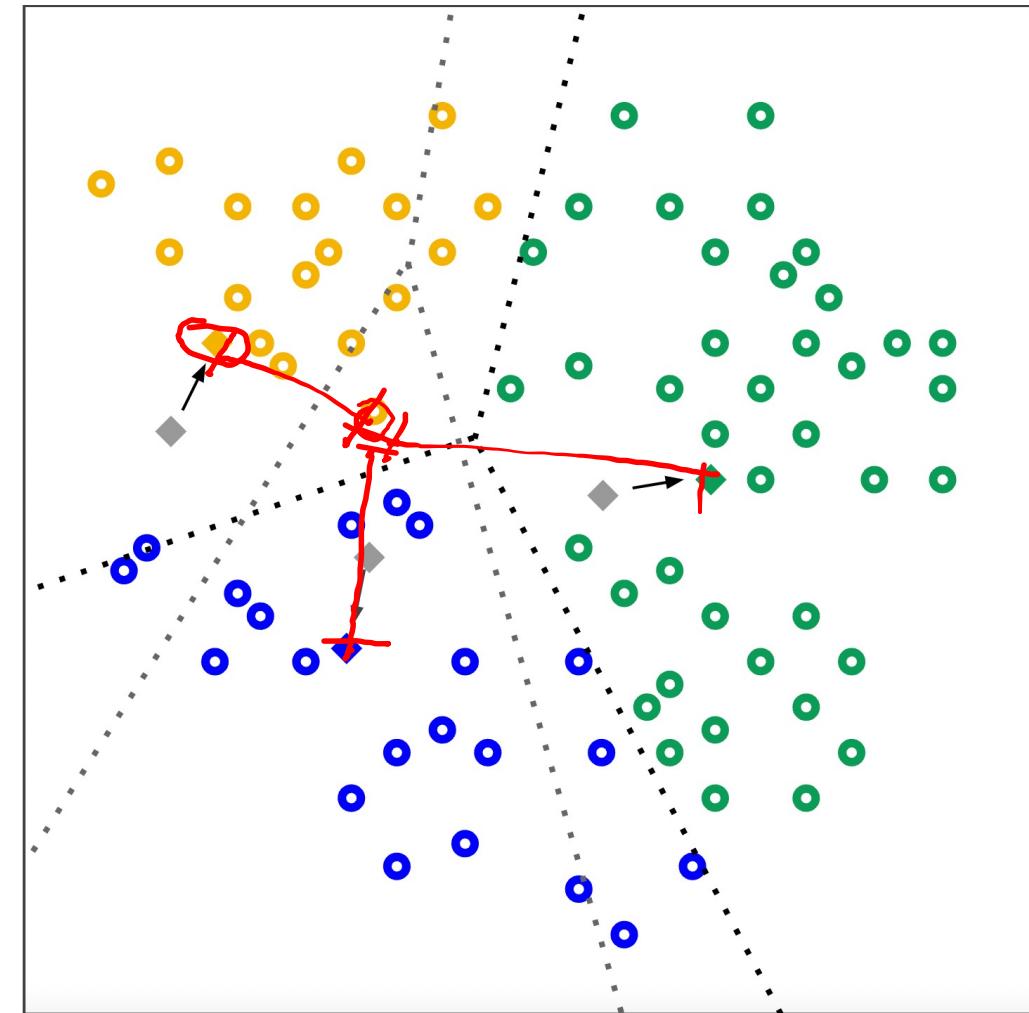
Hitung ulang posisi centroid

- Untuk setiap cluster, hitung ulang posisi centroid dengan cara menghitung rata-rata dari seluruh data pada cluster tersebut
- Langkah ini akan merubah posisi centroid lama ke baru (centroid bergeser dari titik abu-abu ke titik yang baru)



Cluster baru sesuai centroid baru

- Karena centroid telah berubah, maka posisi seluruh data terhadap cluster baru juga berubah
- Anggota cluster dihitung ulang menurut centroid terdekat (pembagian cluster disesuaikan menurut posisi centroid yang baru)

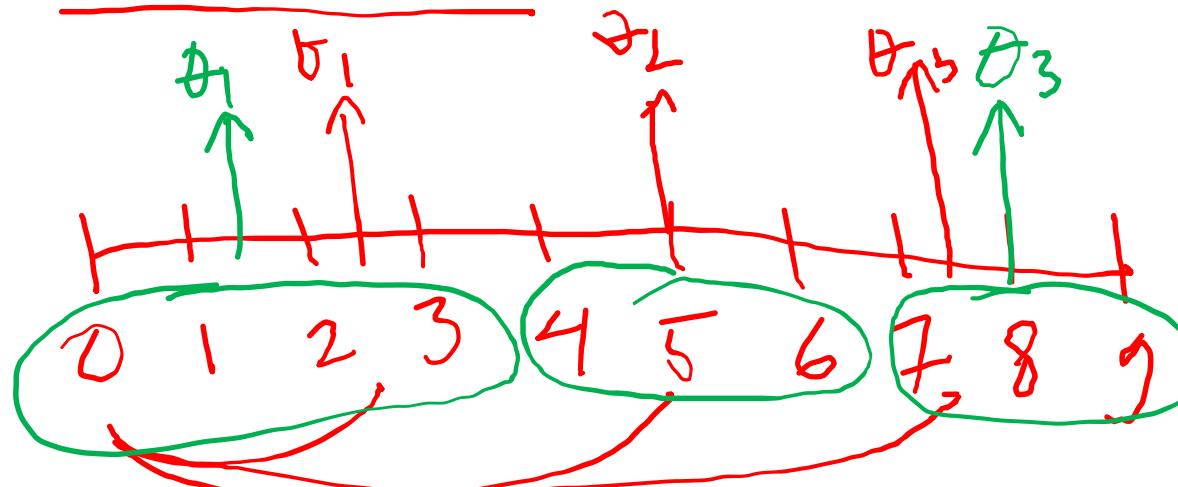


Uji Pemahaman

1 1 1 2 2 2 3 3 3
↑↑↑↑↑↑↑↑↑↑↑↑↑↑↑↑

K=3

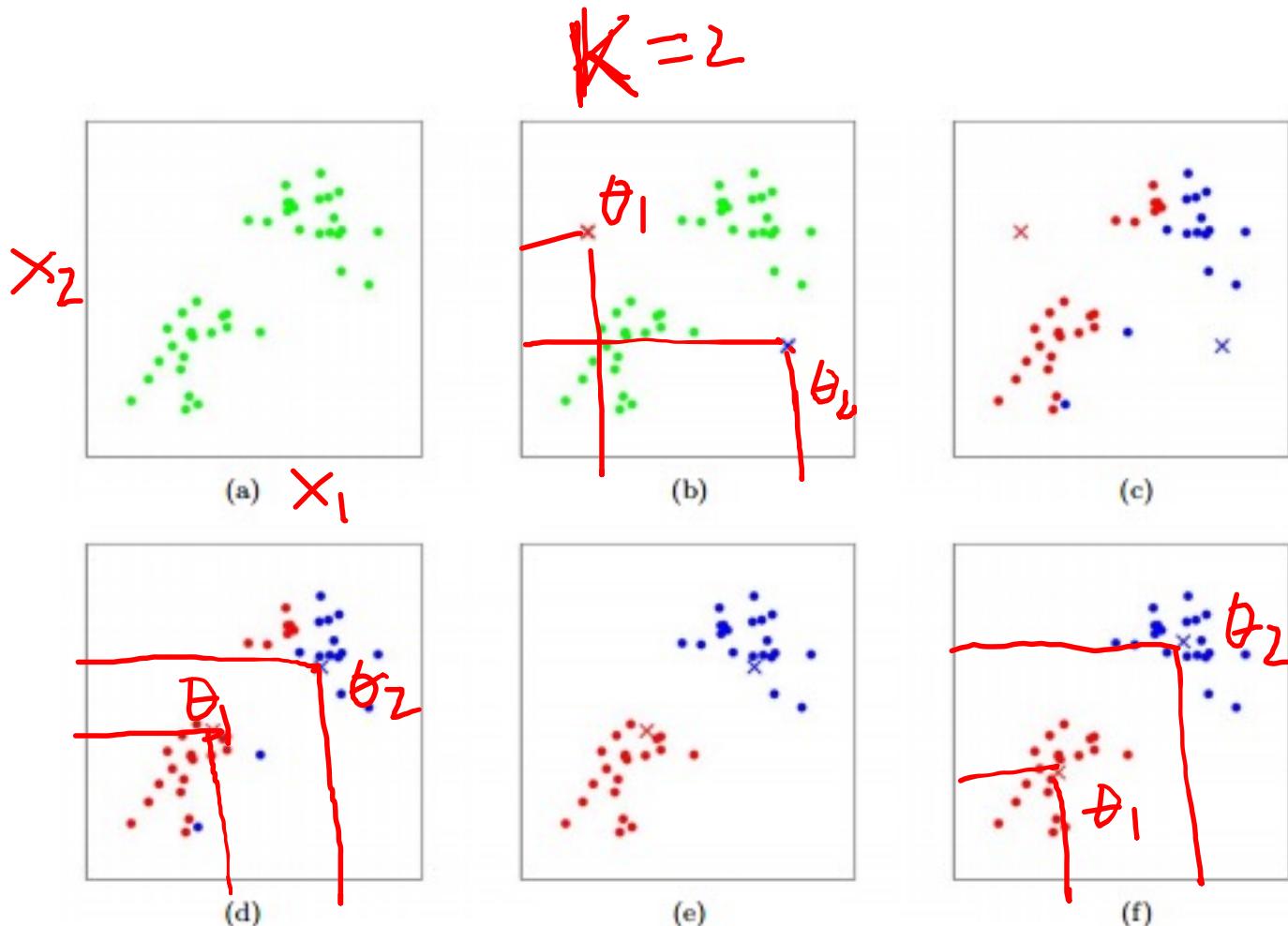
- Terdapat data 1 fitur $X = [0,1,2,3,4,5,6,7,8,9]$
- Untuk k-means clustering 3 cluster dengan centroid ($\theta_1 = 0.25$, $\theta_2 = 0.5$, $\theta_3 = 0.75$), manakah yang termasuk cluster 1, 2, dan 3?
- Berapakah nilai centroid yang baru untuk arrangement cluster tersebut?



$$\theta_1 = \frac{0+1+2+3}{4} = 1.5 \quad \theta_2 = \frac{4+5+6}{3} = 5 \quad \theta_3 = \frac{7+8+9}{3} = 8$$

Iterasi K-means

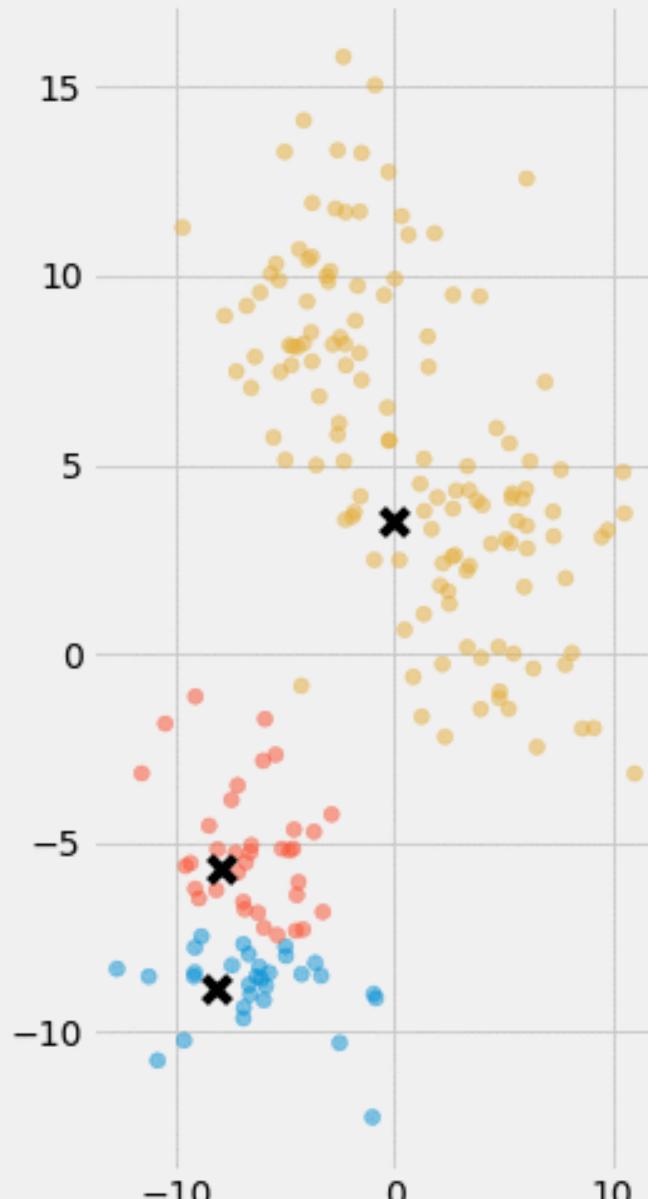
- <https://stanford.edu/~cziegler/cs221/handouts/kmeans.html>
- (a) Original dataset. (b) Random initial cluster centroids. (c-f) Illustration of running two iterations of k-means. In each iteration, we assign each training example to the closest cluster centroid
- Images courtesy of Michael Jordan.



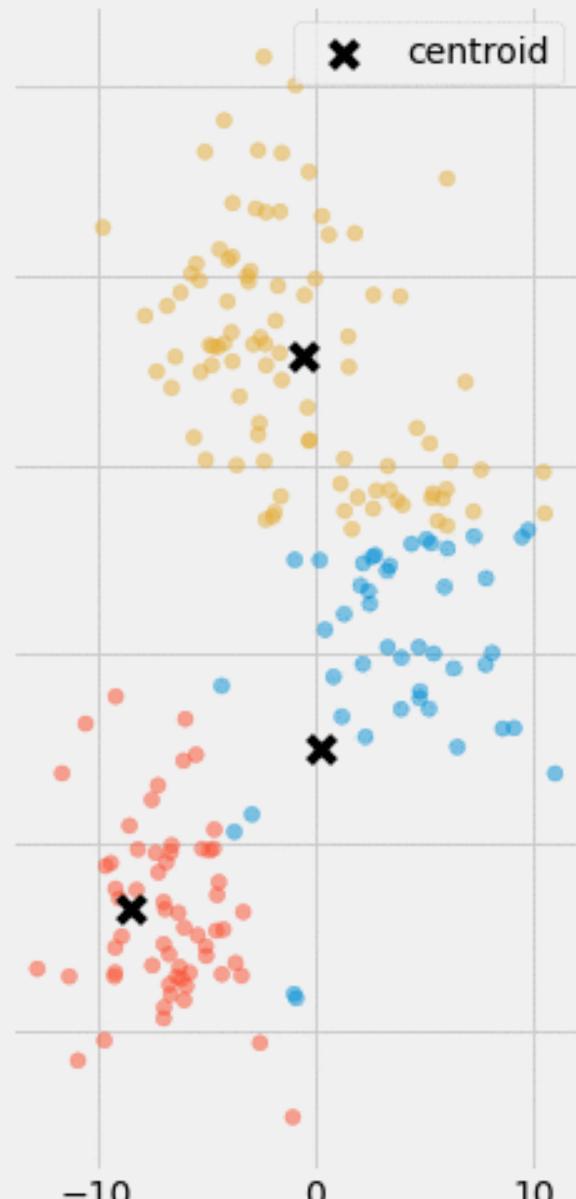
K-means konvergen

- K-means sensitive terhadap posisi awal centroids
- K-means tidak selalu konvergen pada hasil yang sama

k-means iteration: 1
Initialization #1
SSE: 6735.8



Initialization #2
SSE: 5238.5

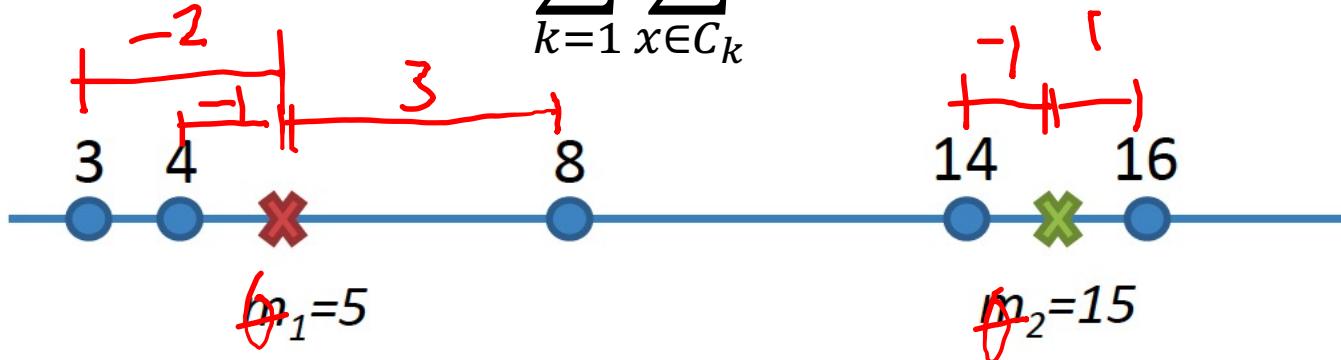


Evaluasi

- Sum of square error (SSE)
 - Untuk setiap titik data, error (jarak dari cluster terdekat) dihitung
 - Semua error dikuadratkan lalu dijumlah

$$SSE = \sum_{k=1}^K \sum_{x \in C_k} (\cancel{m}_k - x)^2$$

- Contoh:



$$SSE = (3 - 5)^2 + (4 - 5)^2 + (8 - 5)^2 + (14 - 15)^2 + (16 - 15)^2$$

$$SSE = (4 + 1 + 9) + (1 + 1) = 16$$

sse

Uji Pemahaman

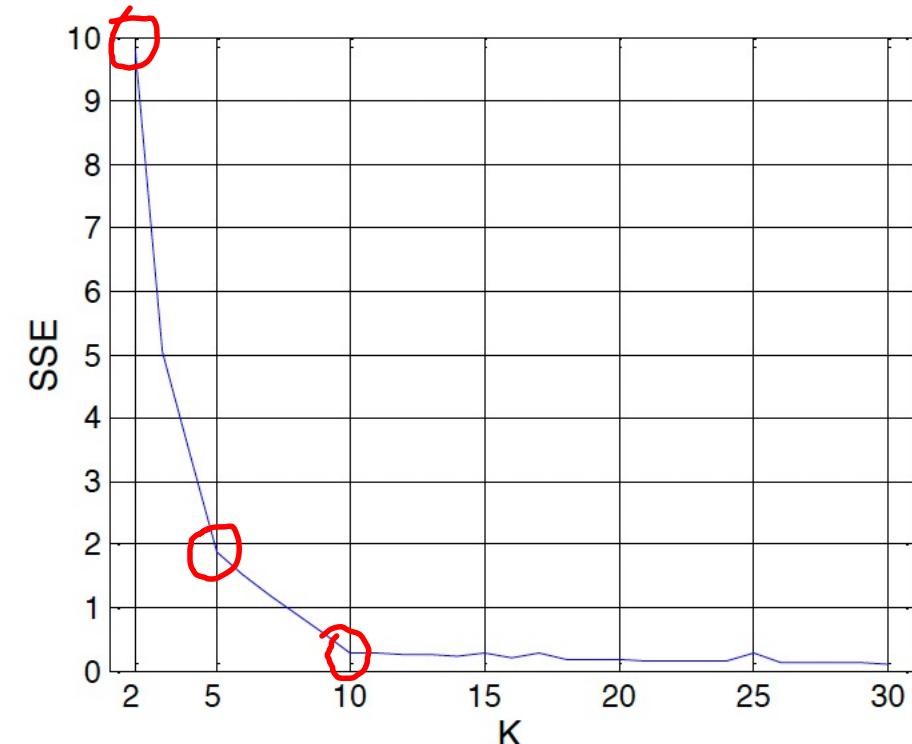
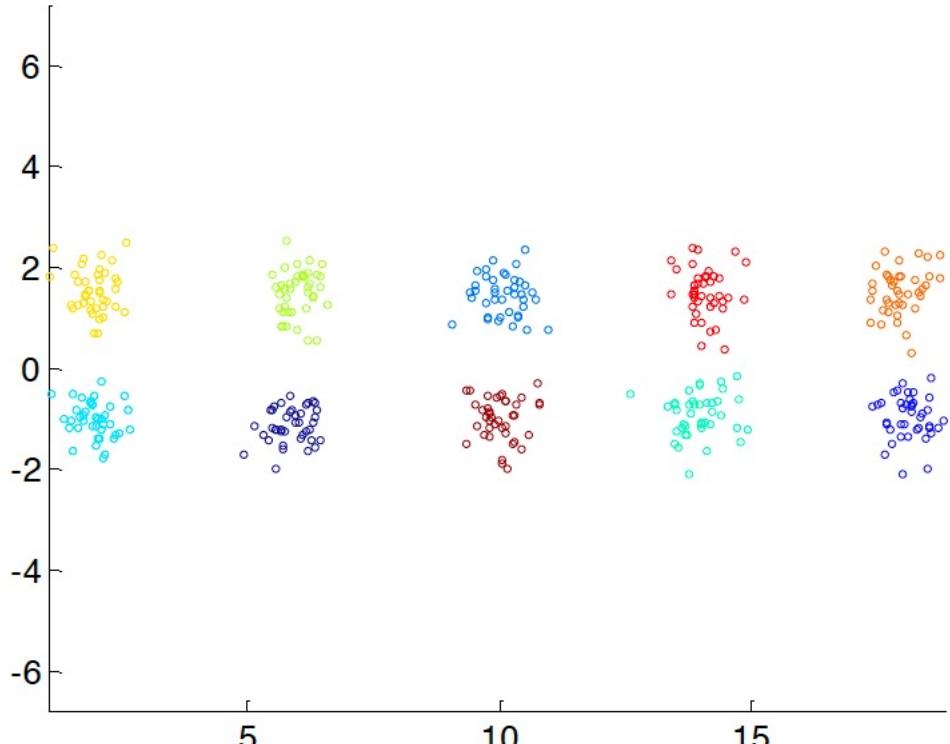
- Berapa nilai SSE dari data 1 fitur $X = [0,1,2,3,4,5,6,7,8,9]$ dengan centroid ($\mu_1 = 2.5$, $\mu_2 = 5$, $\mu_3 = 7.5$)?

$$(0-2.5)^2 + (1-2.5)^2 + (2-2.5)^2 + (3-2.5)^2 + (4-5)^2 + (5-5)^2 + (6-5)^2 + (7-7.5)^2 + (8-7.5)^2 + (9-7.5)^2$$

$$6.25 + 2.25 + 0.25 + 0.25 + 1 + 0 + 1 + 0.25 + 0.25 + 2.25 \\ \underbrace{4 \times 0.25}_{1} + \underbrace{2 \times 1}_{2} + \underbrace{2 \times 2.25}_{4.5} + 6.25 = 13.75$$

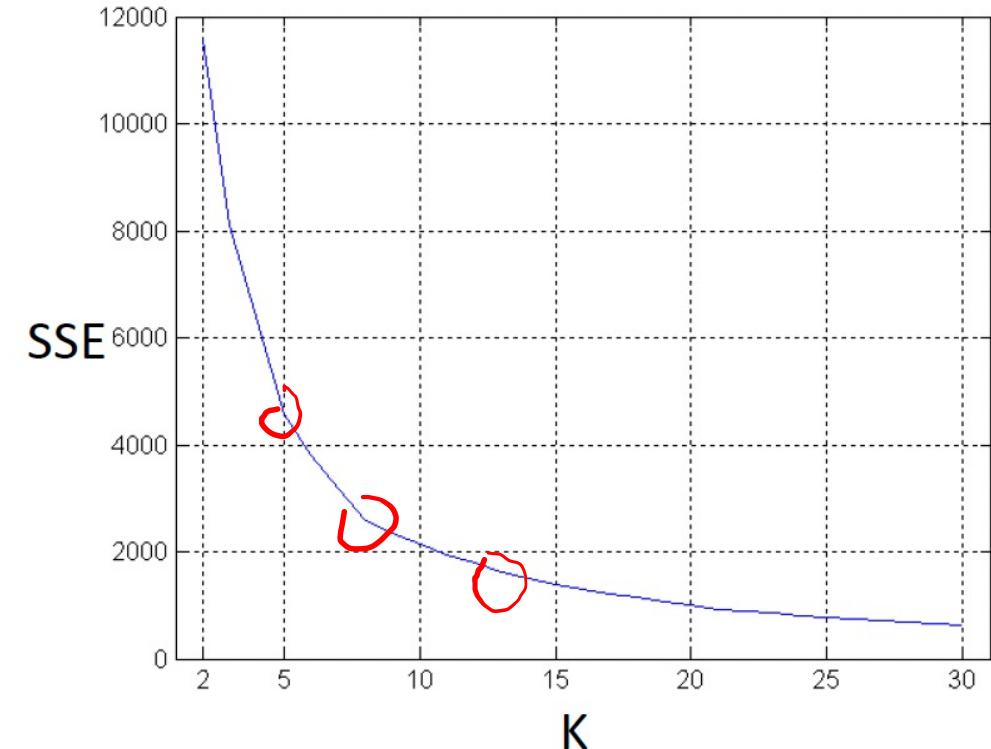
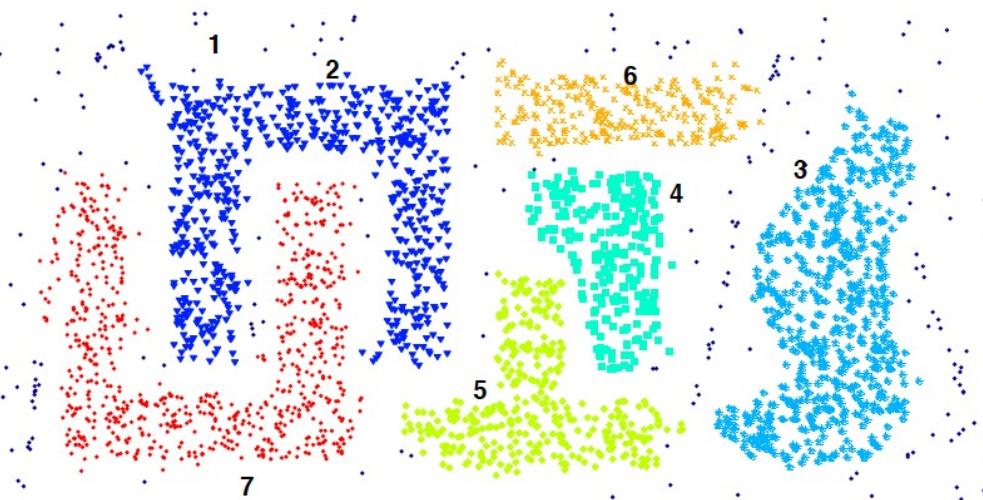
Menentukan K optimal

- Plot SSE untuk setiap pilihan jumlah cluster K
- Elbow method (pilih nilai K pada siku kurva)



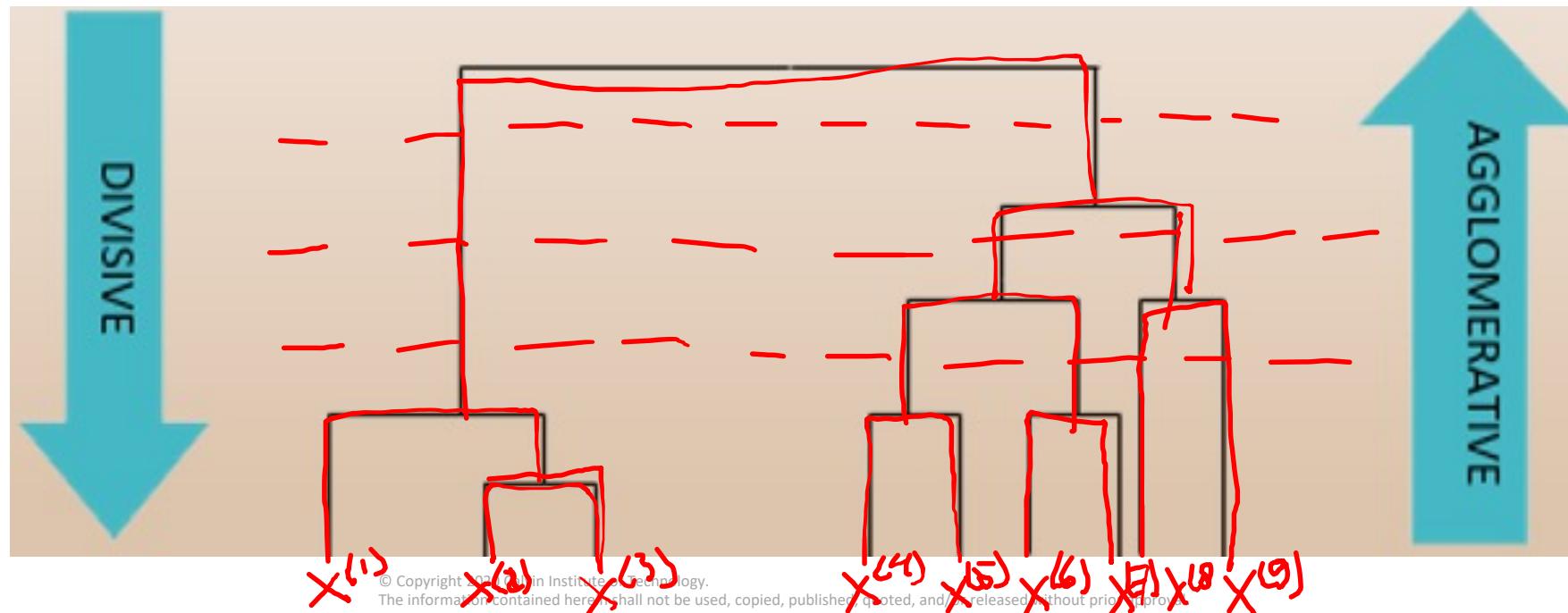
Nilai K ambigu

- Terkadang sulit menentukan siku kurva SSE
- Gunakan metode lain seperti metode silhouette
- Gunakan algoritma lain seperti DBScan, hierarchical clustering, spectral clustering, dll



Hierarchical Clustering

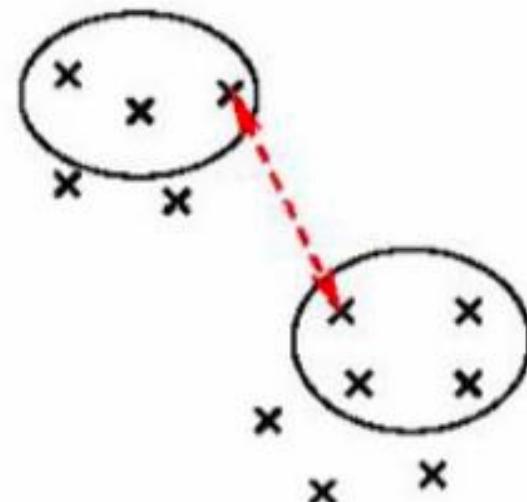
- Membangun hirarki cluster dari tunggal menjadi banyak atau sebaliknya
- Terdapat 2 jenis strategi pengelompokan:
 - Agglomerative: dimulai dengan setiap data sebagai sebuah cluster kemudian membentuk cluster yang semakin membesar
 - Divisive: dimulai dengan semua data dikumpulkan dalam satu cluster kemudian dibagi sampai setiap data berada pada cluster tersendiri



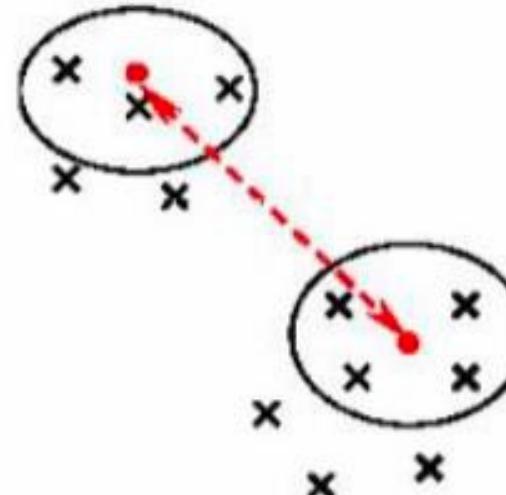
Agglomerative Hierarchical Clustering

- Teknik pengelompokan:
 - Single linkage (jarak terdekat): menggabungkan cluster menurut jarak 2 data terdekat antar cluster
 - Average linkage (jarak rata-rata): menggabungkan cluster menurut jarak rata-rata tiap pasangan data antar cluster
 - Complete linkage (jarak terjauh): menggabungkan cluster menurut jarak 2 data terjauh antar cluster

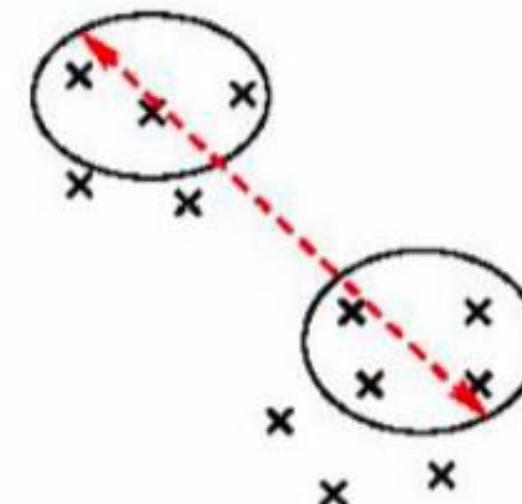
- Simple linkage



- Average linkage



- Complete linkage



Uji Pemahaman

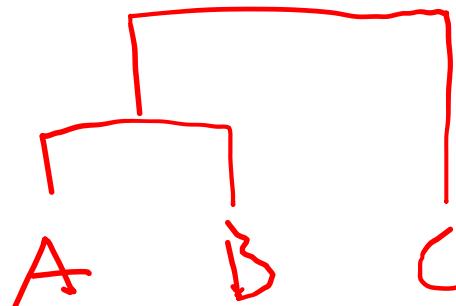
- Terdapat data 1 fitur yang sudah membentuk clusters:

- $A = [0,1,2,3,4]$
- $B = [5,6,7]$
- $C = [8,9]$

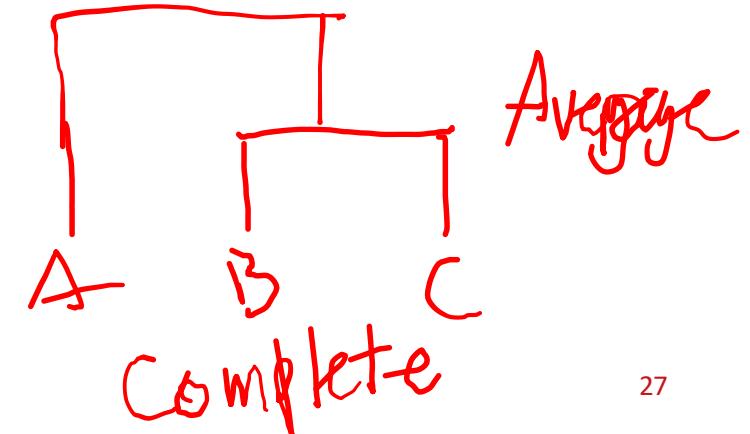
$\{ \begin{matrix} 4 \\ 2,5 \end{matrix} \} \quad \{ \begin{matrix} 7 \\ 4 \end{matrix} \} \quad \} y$

- Manakah bentuk pengelompokan yang tepat menurut:

- Single linkage
- Average linkage
- Complete linkage



Single



Complete

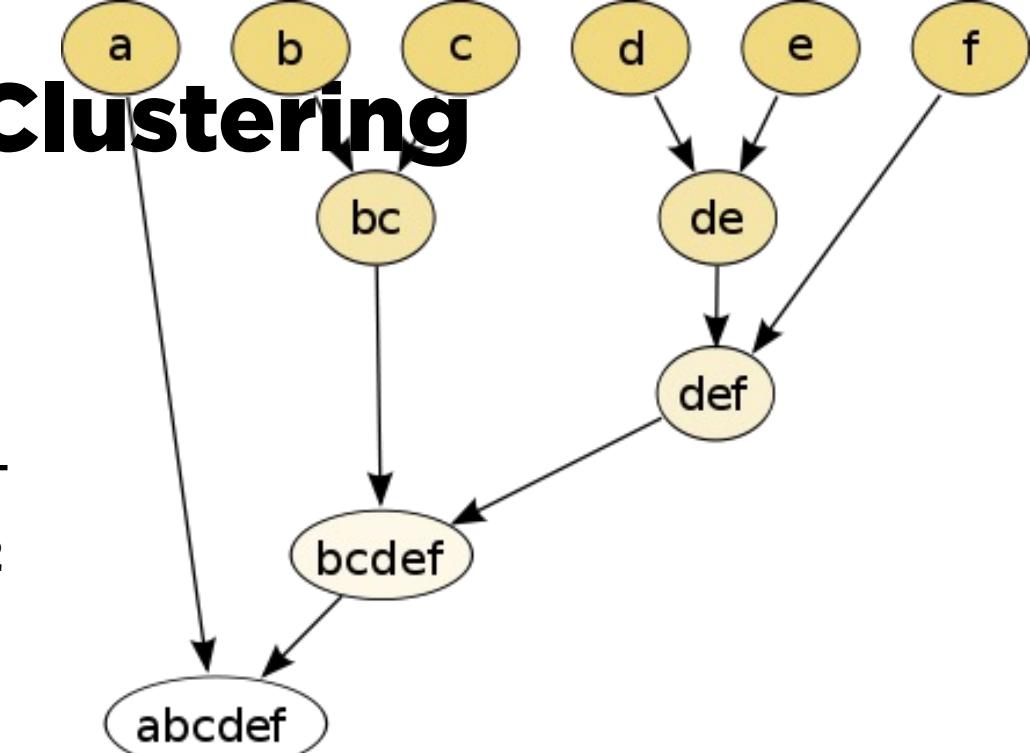
Average

Agglomerative Hierarchical Clustering

- Algoritma:
 - Buat matriks jarak dengan menghitung jarak
 - Misalnya jarak Euclidean:

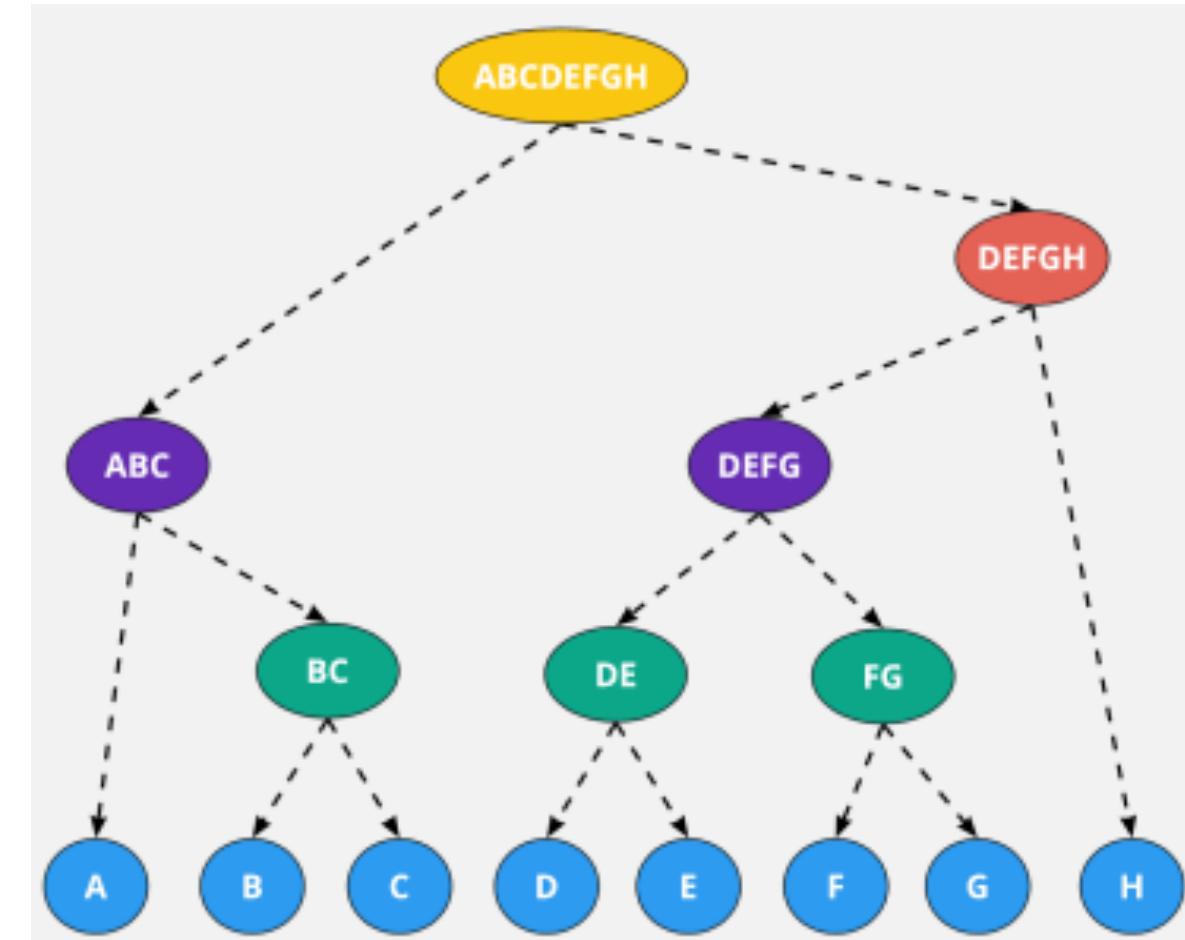
$$d_{ii'} = \sqrt{\sum_{j=1}^N (x_j^{(i)} - x_j^{(i')})^2}$$

- Gabungkan 2 cluster terdekat: jika data a dan b adalah 2 data terdekat maka akan membentuk cluster gabungan pertama
- Perbaharui matriks jarak sesuai dengan teknik pengelompokan (single, average, complete)
- Ulangi Langkah 2-3 hingga tersisa 1 cluster yang memuat seluruh data
- Buat dendrogram (“dendron”=pohon, “gramma”=gambar)



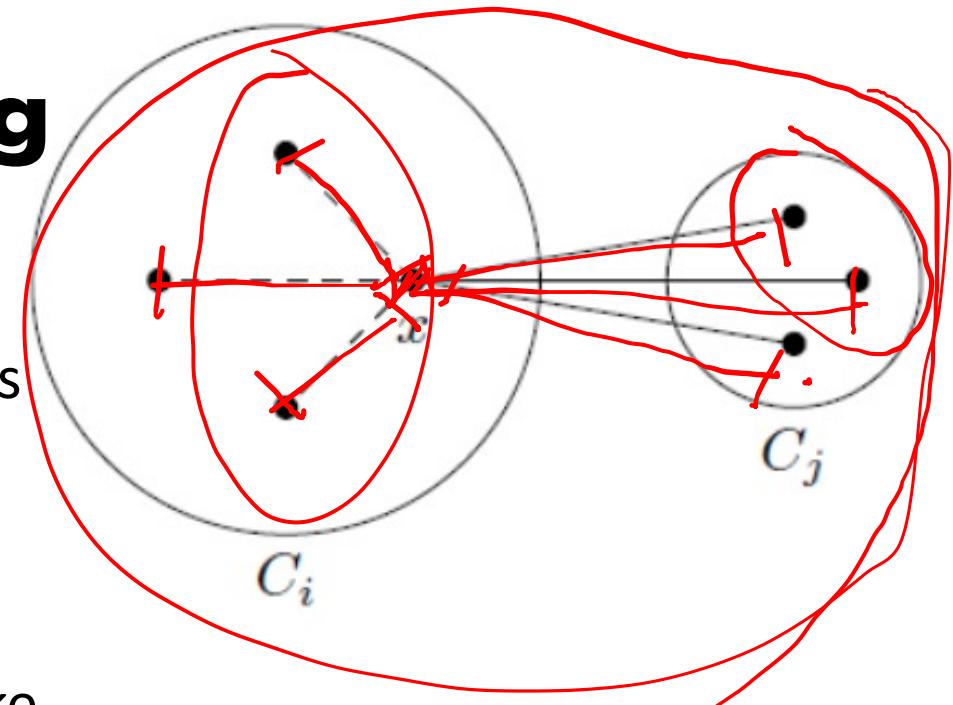
Divisive Hierarchical Clustering

- DIANA (Divisive ANalysis):
 - membagi cluster tunggal menjadi 2 clusters yang lebih kecil, dan seterusnya hingga setiap data adalah 1 cluster tersendiri



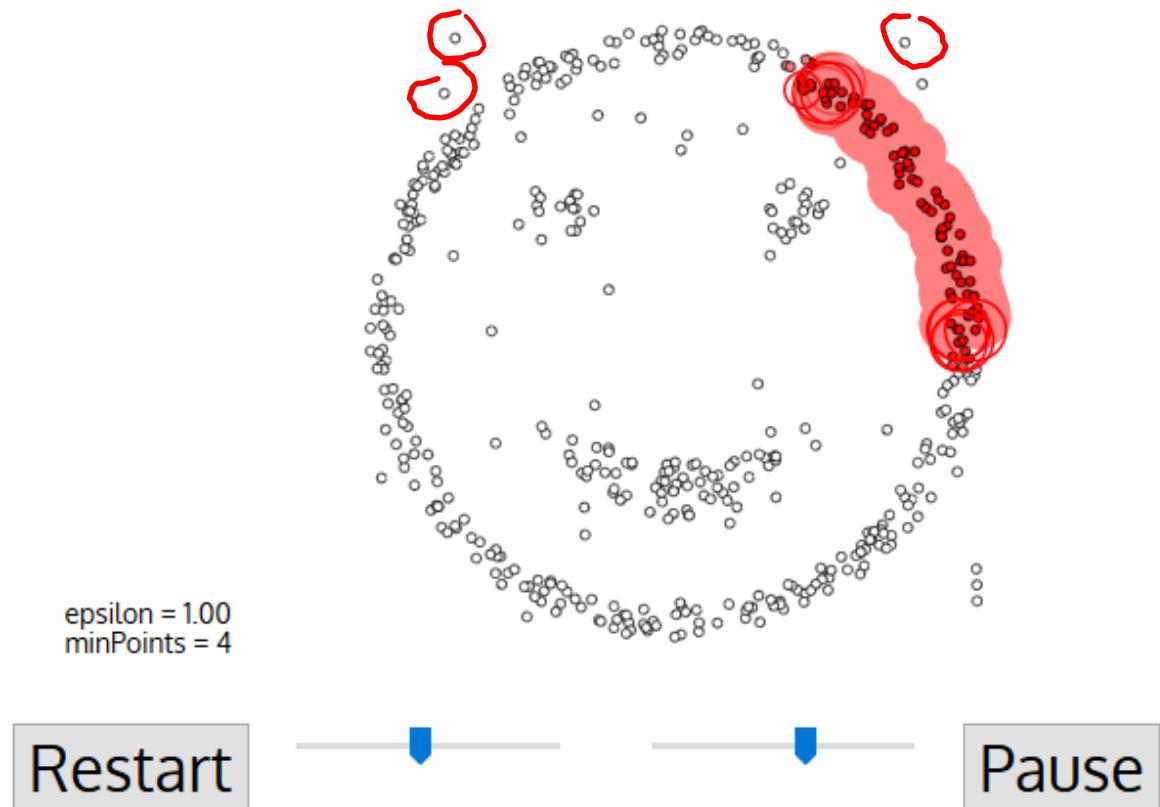
Divisive Hierarchical Clustering

- Algoritma:
 - Suatu cluster tunggal C akan dibagi menjadi clusters C_i dan C_j
 - Anggaplah $C_i = C$ sedangkan $C_j = \emptyset$
 - Untuk setiap data $x \in C_i$:
 - Untuk iterasi pertama: hitung jarak rata-rata \underline{x} ke semua data \rightarrow pindah data dengan jarak rata-rata maksimum ke C_j
 - Untuk iterasi selanjutnya: hitung $D_x = avg\{d(x, y), y \in C_i\} - avg\{d(x, y), y \in C_j\} \rightarrow$ tentukan data di C_i yang memiliki D_x paling besar
 - Ulangi sampai semua nilai D_x negative
 - Tentukan cluster terkecil dengan diameter terbesar.
 - Ulangi hingga setiap cluster hanya ada 1 data



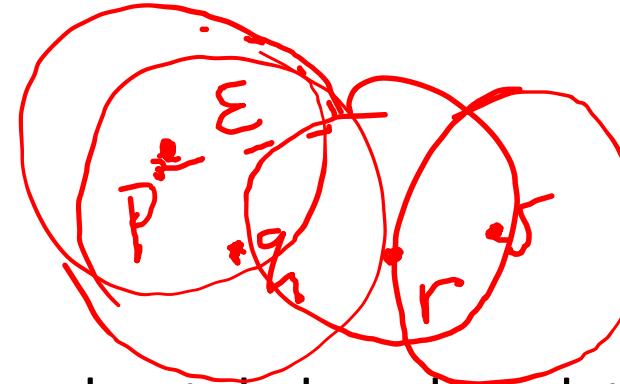
DBScan

- Density Based Spatial Clustering of Applications with Noise
- <https://www.kdnuggets.com/2020/04/dbscan-clustering-algorithm-machine-learning.html>
- Mengelompokan data yang berdekatan satu sama lain (high density)
- Data di daerah low density (tetangga jauh) → outliers / noise
- Cluster: wilayah terhubung dengan kepadatan tinggi
- Bagaimana menentukan anggota cluster?
Epsilon? Minpoints?



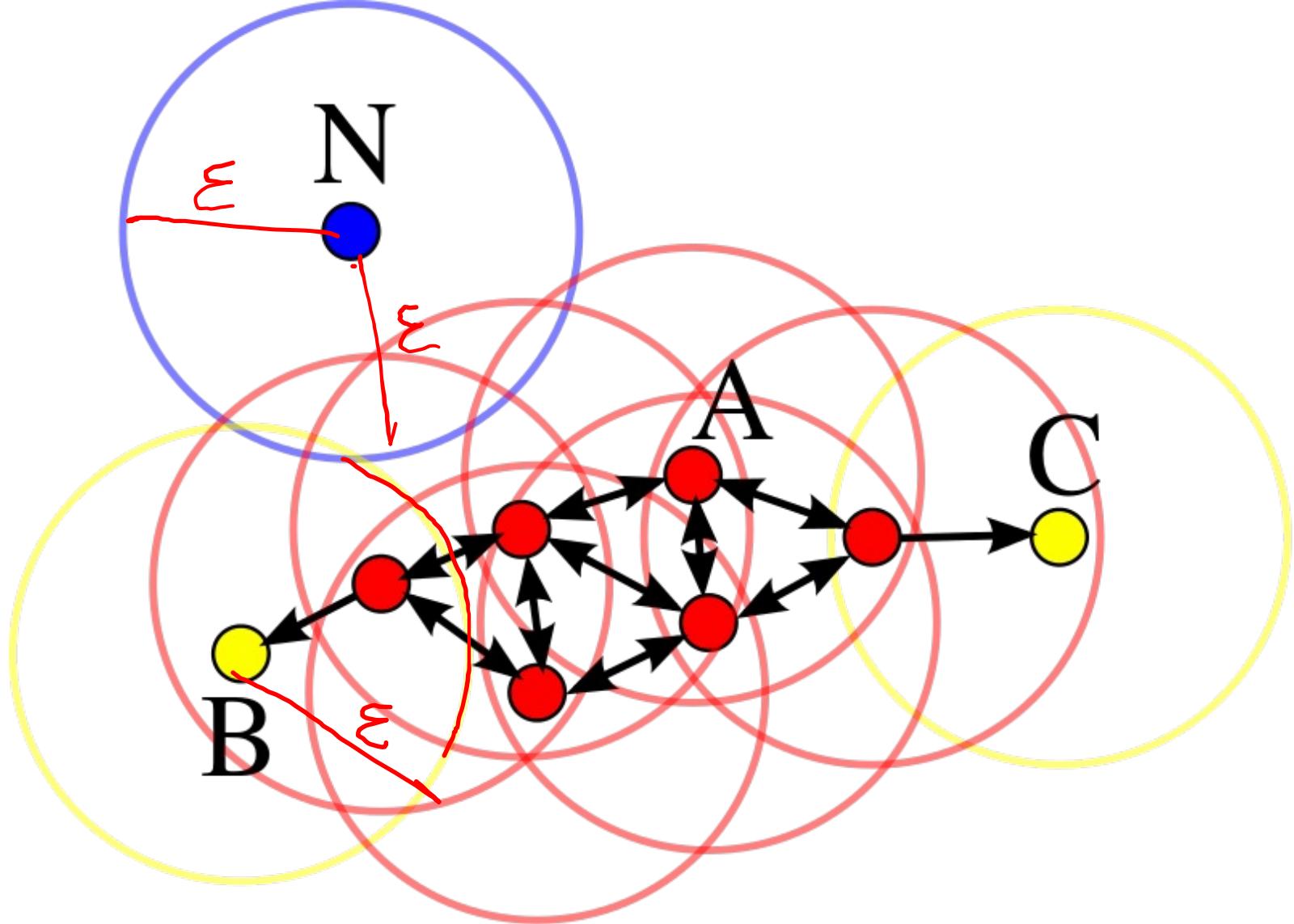
DBScan

- Definisi:
 - ε : radius dari suatu titik p
 - MinPts: jumlah minimal data untuk membentuk daerah padat
- Penggolongan data:
 - Core Point → titik p disebut core point jika ada minimal MinPts titik yang berjarak $\leq \varepsilon$ dari p (termasuk p sendiri) → dense region
 - Directly reachable point → titik q directly reachable jika berjarak $\leq \varepsilon$ dari p
 - Reachable point → titik q reachable dari p jika ada jalur p_1, p_2, \dots, p_n dimana di awal jalur $p_1 = p$, di akhir jalur $p_n = q$, dan di tengah-tengah jalur p_{i+1} directly reachable dari p_i
 - Outliers/noise points → titik yang tidak reachable dari titik lain manapun



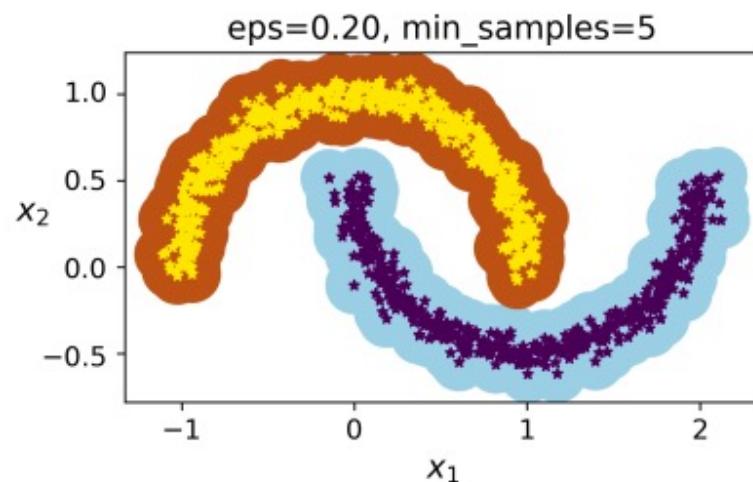
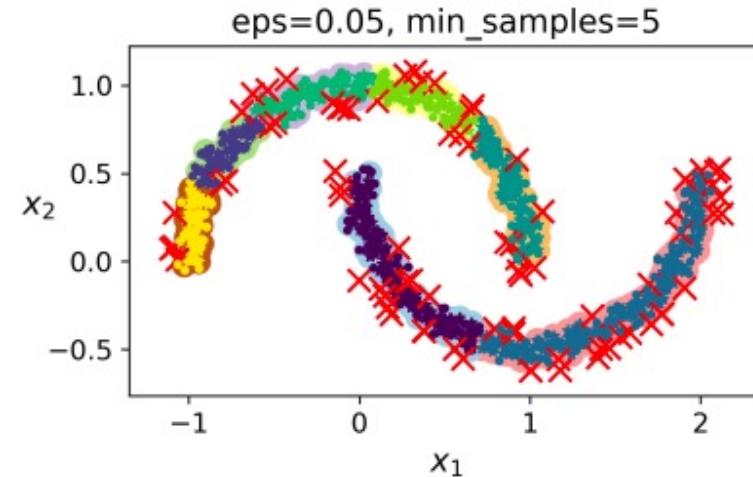
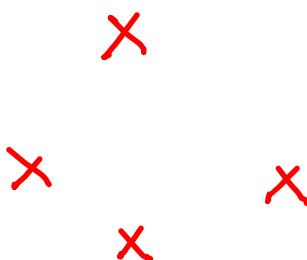
DBScan

- Lingkaran = radius ε
- MinPts = 4
- Merah = core points
- Kuning = reachable points
- Biru = outliers / noise points



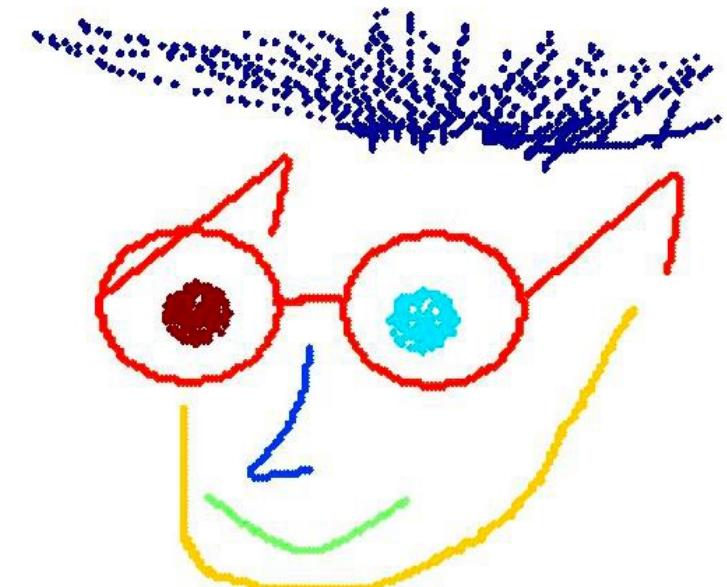
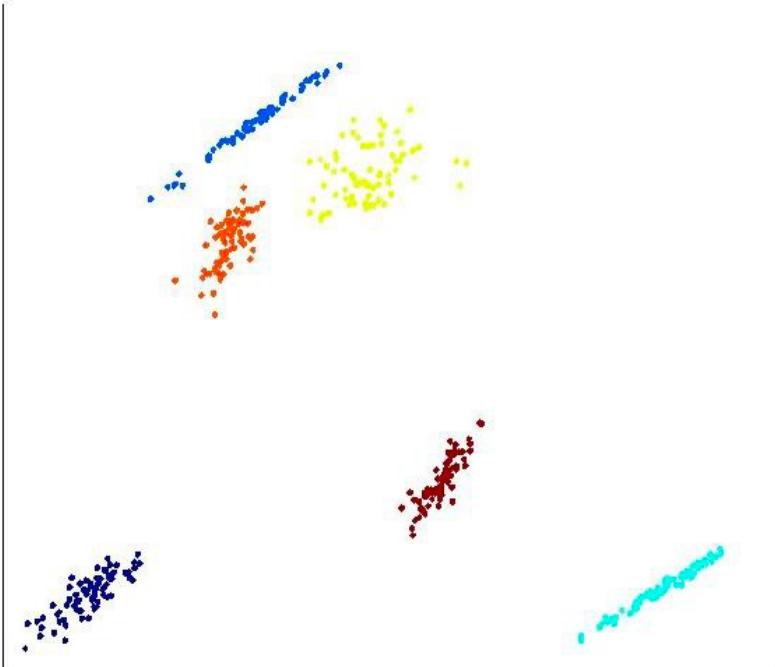
DBScan

- Algoritma bekerja baik jika semua clusters cukup padat dan dipisahkan oleh wilayah low density secara jelas
- Kelebihan:
 - Implementasi sederhana
 - Identifikasi clusters berapapun dan bentuk apapun
 - Robust terhadap outliers
 - Hyperparameter hanya 2 (ϵ , minPts)
- Limitasi:
 - Jika kepadatan bervariasi antar clusters
 - Sulit pilih hyperparameter



Spectral Clustering

- Clustering biasa dapat dipakai tetapi model parametrik dan jarak Euclidean tidak selalu cocok



Spectral Clustering

- Affinity Matrix:

$$A_{ii'} = e^{-\frac{\|x^{(i)} - x^{(i')}\|^2}{2\sigma^2}}$$

- Diagonal Matrix:

$$D_{ii} = \sum_{i'} A_{ii'}$$

- Graph Laplacian Matrix:

$$L = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$$

- Lakukan eigen-decomposition terhadap Graph Laplacian Matrix dan kumpulkan k largest eigenvectors:

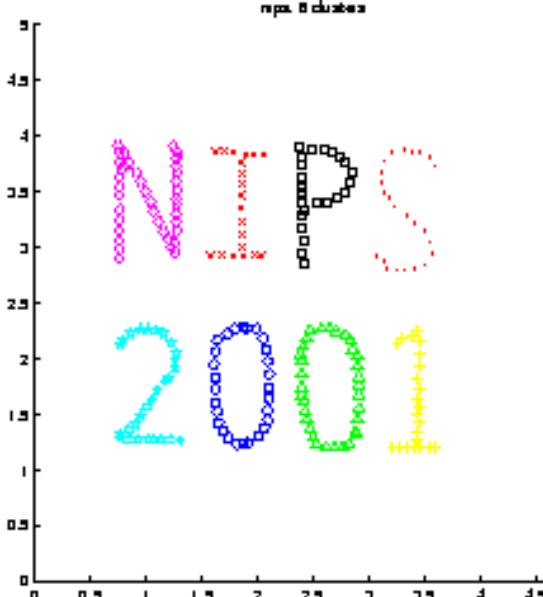
$$v_1, v_2, \dots, v_k$$

- Renormalisasi eigenvector:

$$Y_{ij} = \frac{V_{ij}}{(\sum_j V_{ij}^2)^{\frac{1}{2}}}$$

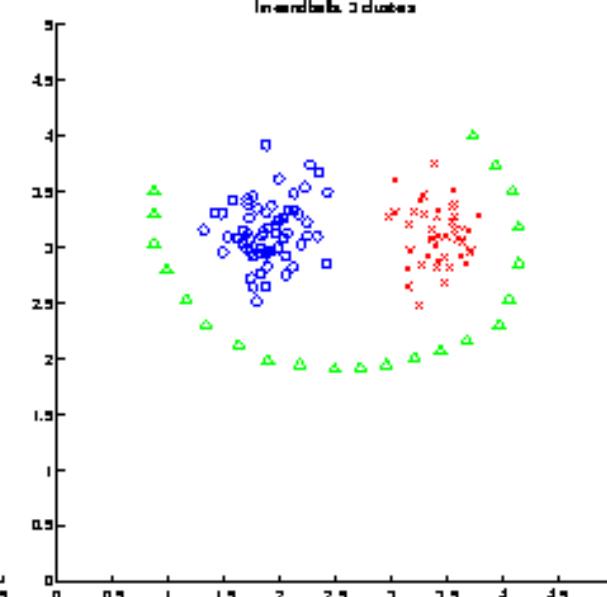
- Cluster $Y \in \mathbb{R}^{N \times k}$ dengan menggunakan K-means

npes. 8 clusters



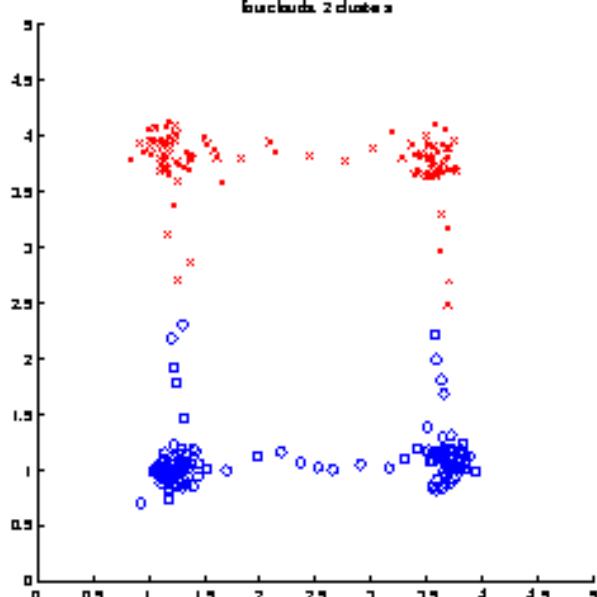
(a)

In-sandels. 3 clusters



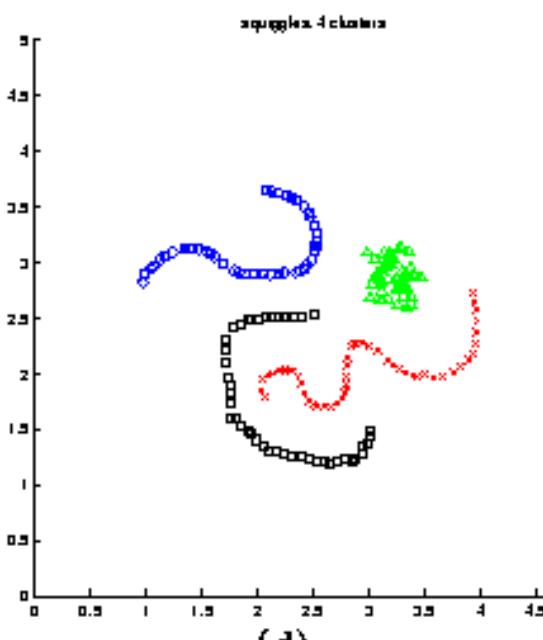
(b)

Guckels. 2 clusters



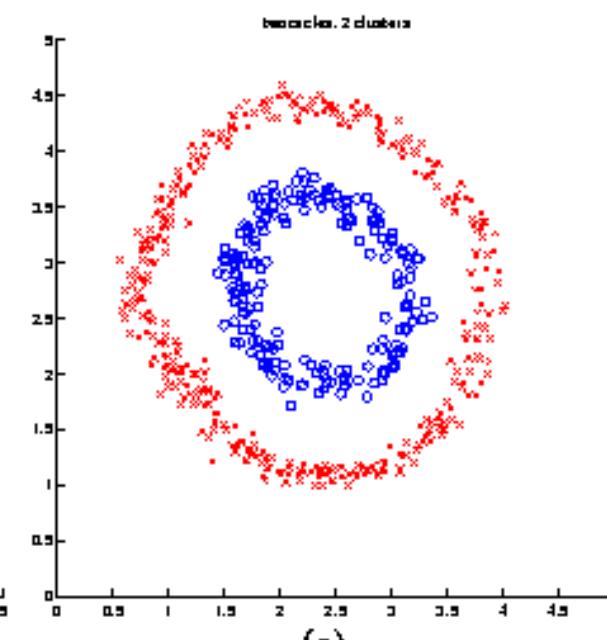
(c)

squiggle. 4 clusters



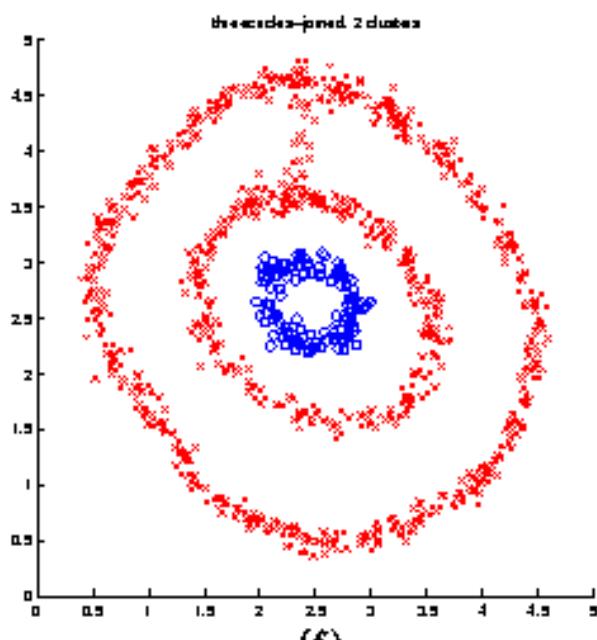
(d)

beeswax. 2 clusters



(e)

three-scales-jointed. 2 clusters



(f)

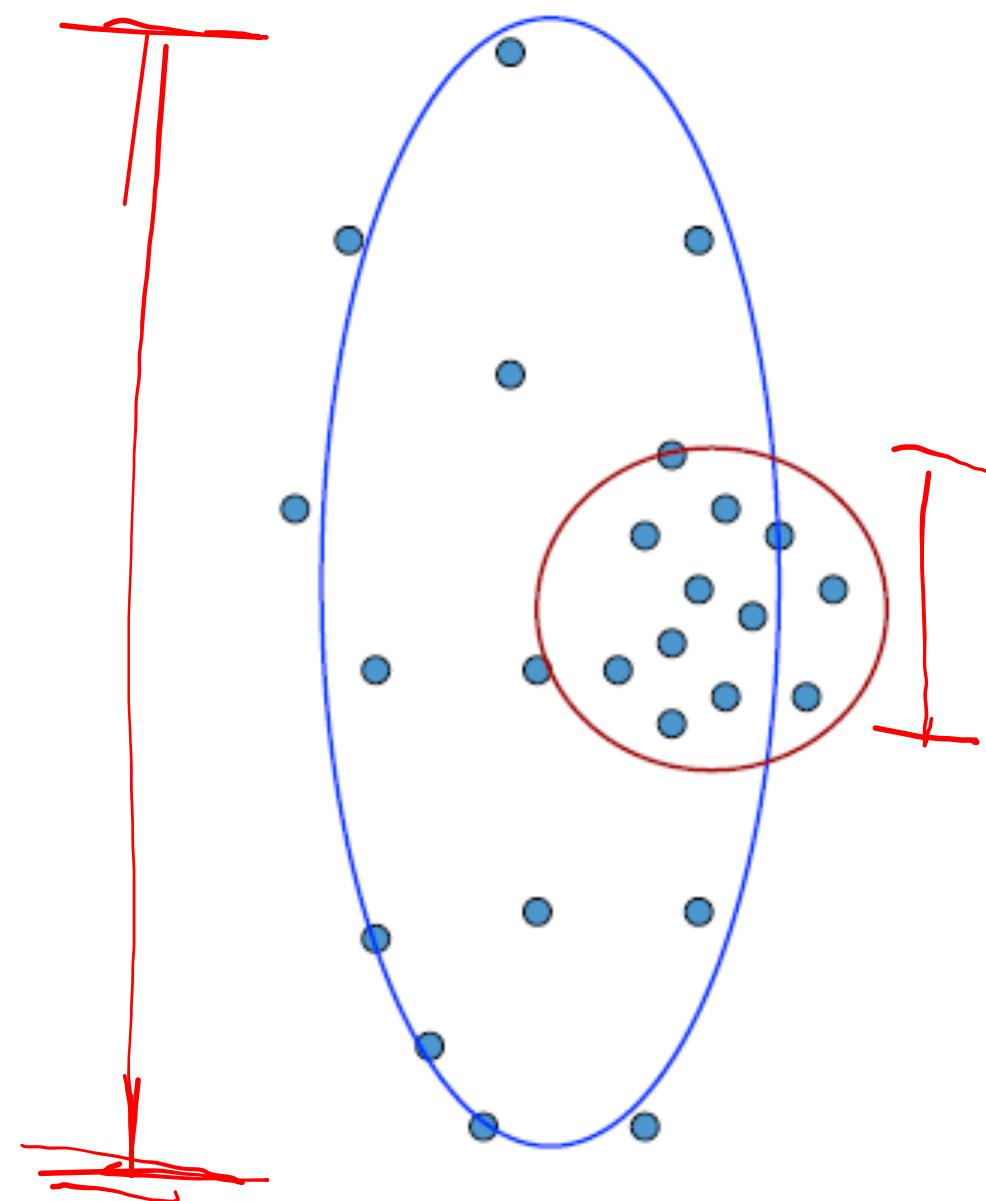
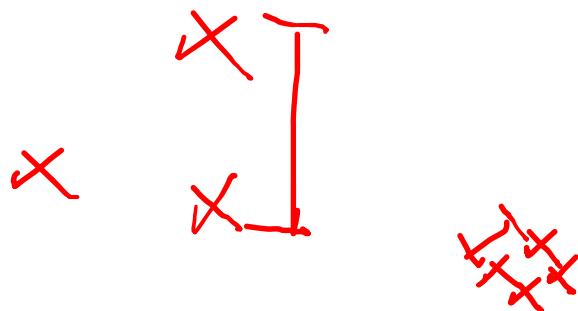
treesca-de-jointed. 3 clusters

Plots of Y (labelled, randomly subsampled) for beeswax

Soft Clustering

Kelemahan Hard Clustering

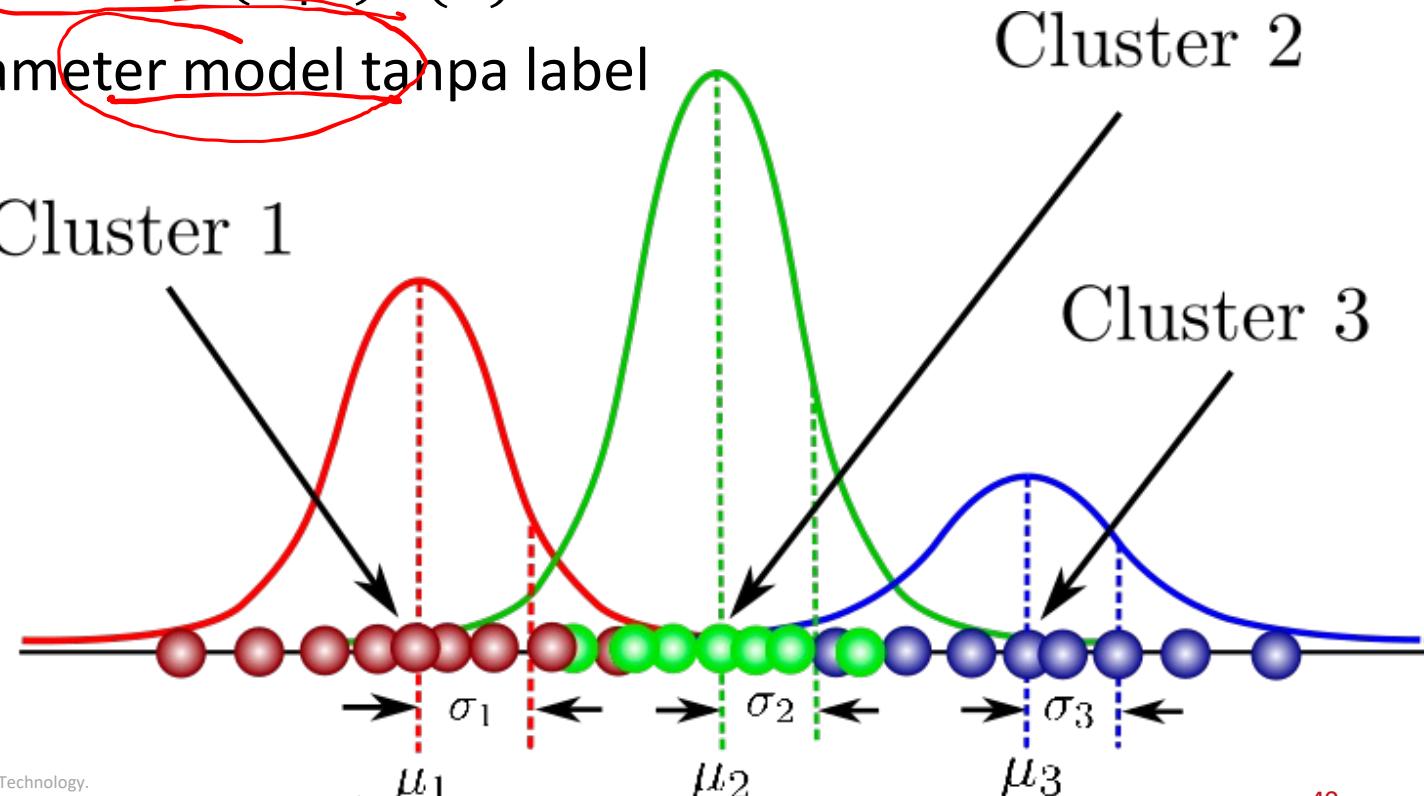
- Cluster bisa overlap
- Beberapa cluster lebih lebar dari yang lain
- Jarak dapat menipu



Probabilistic Clustering

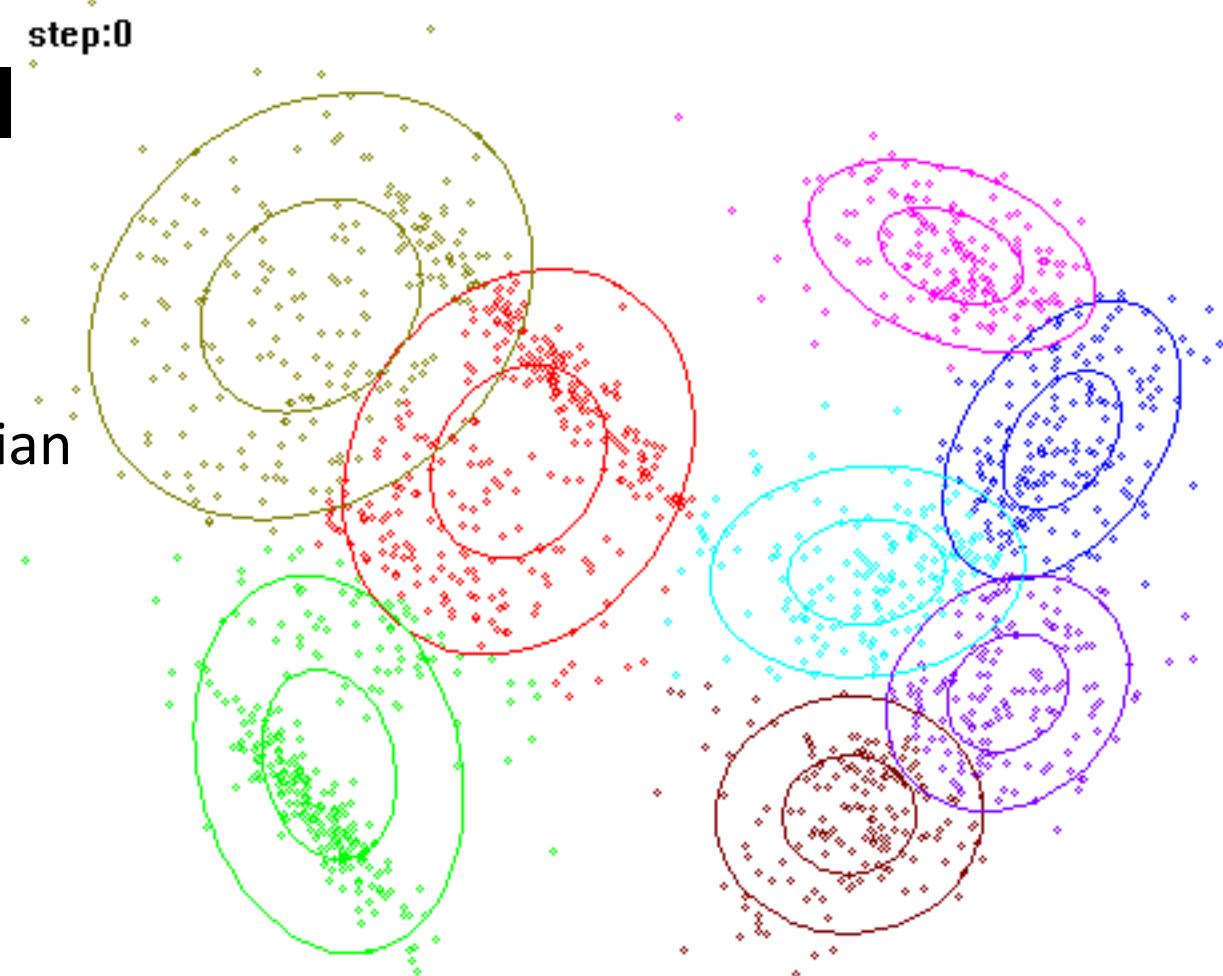
gabungan

- Asumsi: setiap data dapat dibentuk dari mixture beberapa clusters dengan distribusi gaussian yang parameternya dapat dipelajari
- Mengijinkan overlaps clusters berbeda ukuran
- Dapat digunakan untuk generative model: $P(X|Y)P(Y)$
- Challenge: bagaimana estimasi parameter model tanpa label
- Contoh: GMM, PSLI, LDA



Gaussian Mixture Model

- $P(Y)$: terdapat k komponen
- $P(X|Y)$: masing-masing komponen menghasilkan data dari multivariate gaussian dengan mean μ_i dan covariance Σ_i

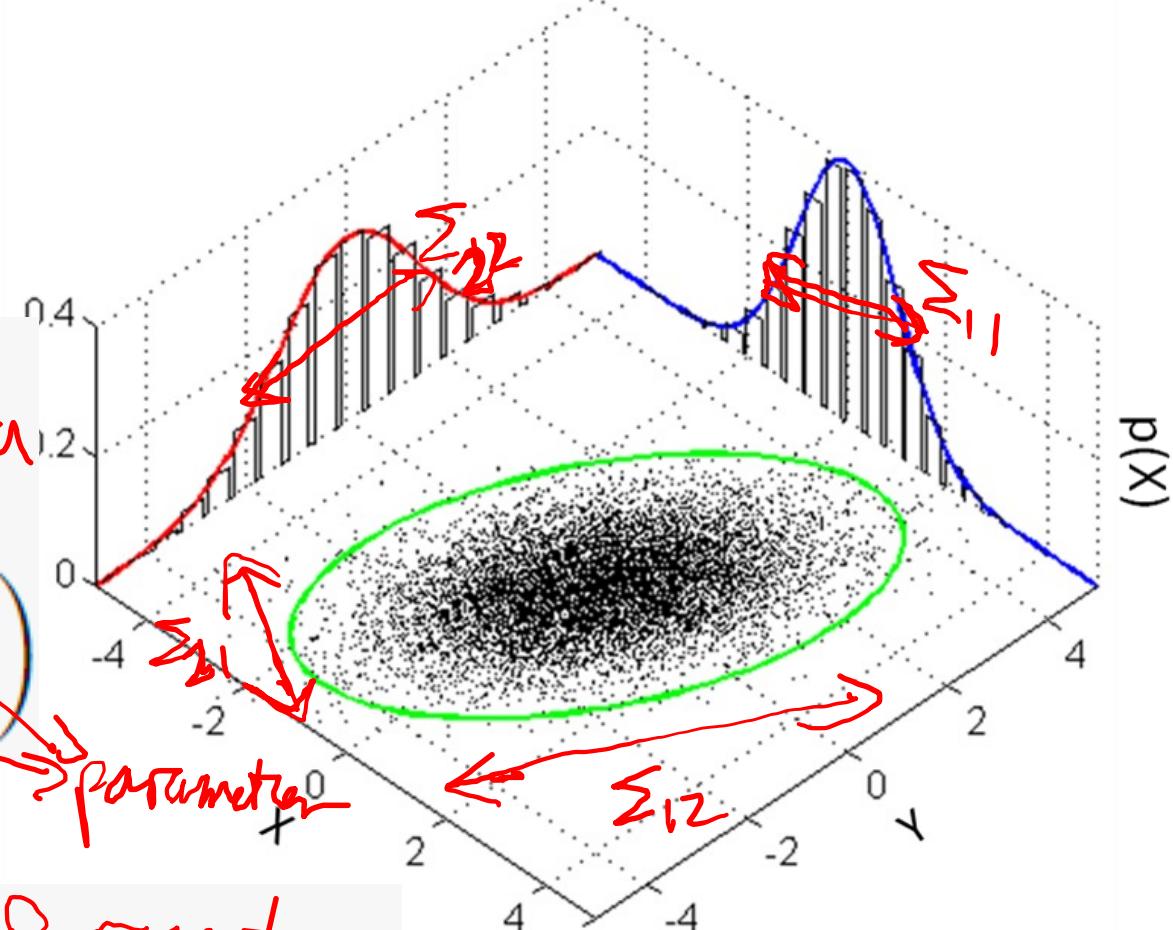


Multivariate Gaussian

- Kasus 1 dimensi

$$p(x) = \sum_{i=1}^K \phi_i \mathcal{N}(x | \mu_i, \sigma_i)$$

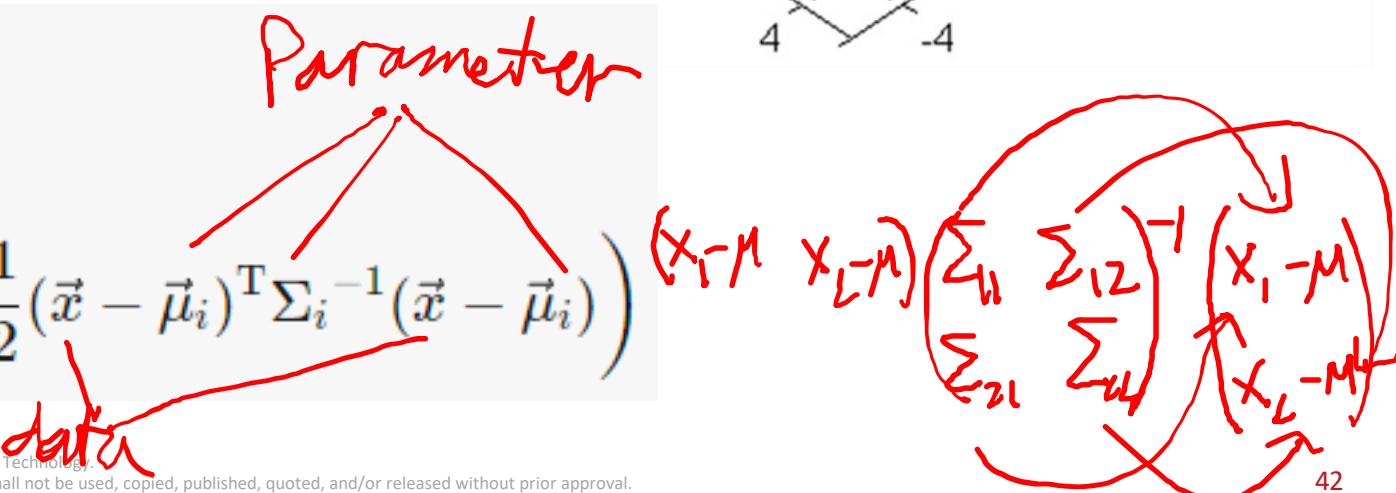
$$\mathcal{N}(x | \mu_i, \sigma_i) = \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left(-\frac{(x - \mu_i)^2}{2\sigma_i^2} \right)$$



- Kasus multi dimensi

$$p(\vec{x}) = \sum_{i=1}^K \phi_i \mathcal{N}(\vec{x} | \vec{\mu}_i, \Sigma_i)$$

$$\mathcal{N}(\vec{x} | \vec{\mu}_i, \Sigma_i) = \frac{1}{\sqrt{(2\pi)^K |\Sigma_i|}} \exp \left(-\frac{1}{2} (\vec{x} - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{x} - \vec{\mu}_i) \right)$$

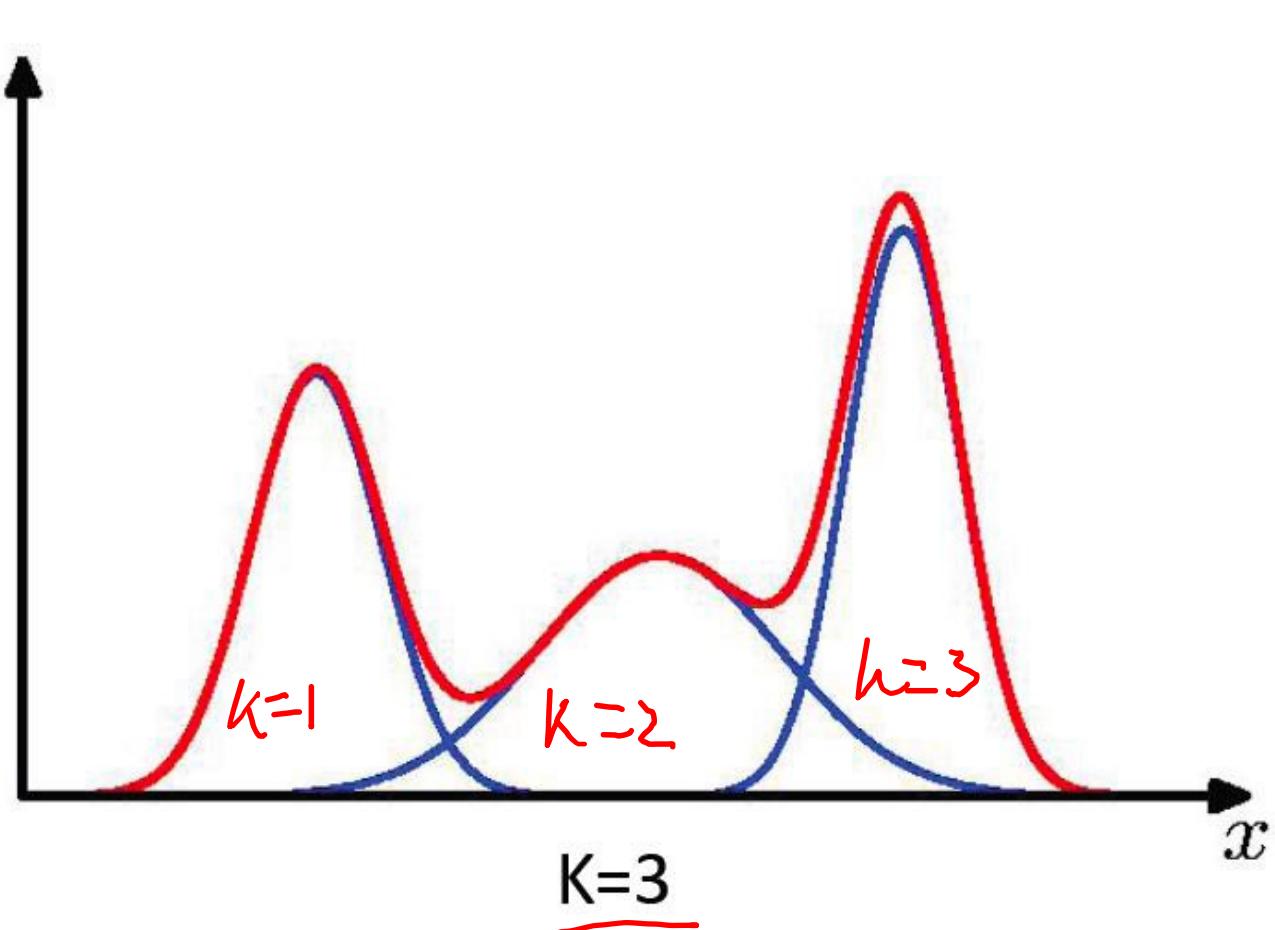


Mixture of Gaussian

- Menjumlahkan beberapa gaussian

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

↑ Component
↑ Mixing coefficient



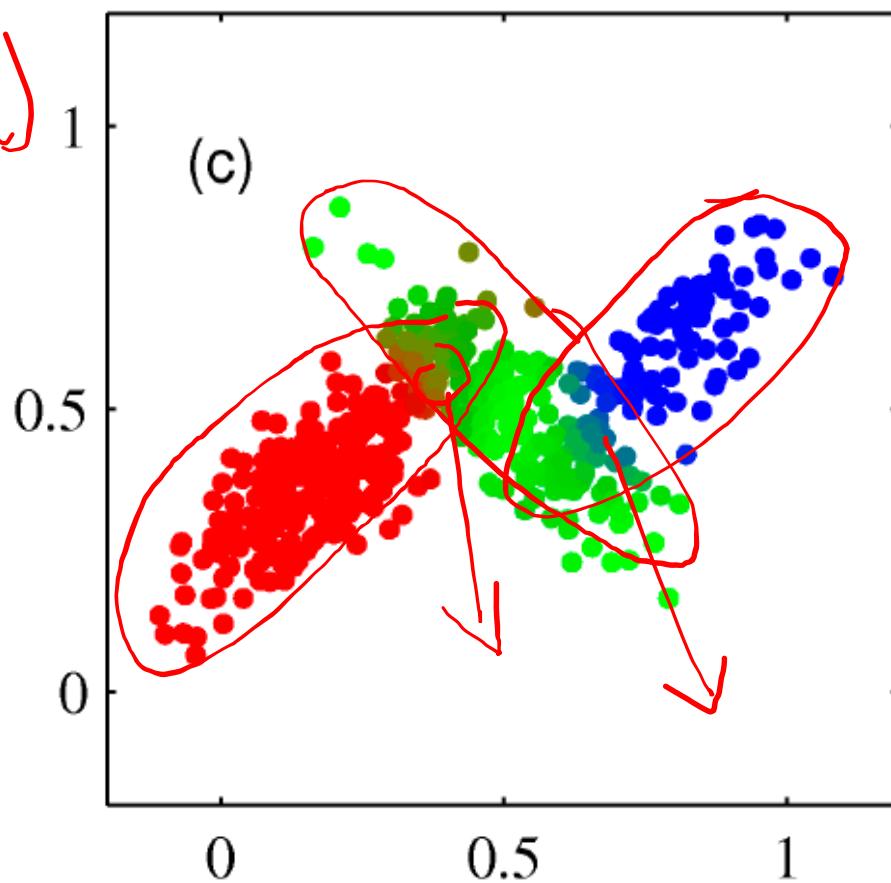
$$\forall k : \pi_k \geqslant 0 \quad \sum_{k=1}^K \pi_k = 1$$

Soft Clustering

- Setiap data merupakan campuran (mixture) dari beberapa multivariate gaussians

$$p(x) = \sum_{k=1}^3 \pi_k N(x | \mu_k, \Sigma_k)$$

$$\sum_{k=1}^3 \pi_k = 1$$



$$D=2 \quad \pi^r, \pi^g, \pi^b$$
$$K=3$$

$$\begin{array}{|c|c|c|c|} \hline & M_1^r, M_1^g, M_1^b \\ \hline & \mu_1^r, \mu_1^g, \mu_1^b \\ \hline & \Sigma_{11}^r, \Sigma_{12}^r, \Sigma_{21}^r, \Sigma_{22}^r \\ \hline & \Sigma_{11}^g, \Sigma_{12}^g, \Sigma_{21}^g, \Sigma_{22}^g \\ \hline & \Sigma_{11}^b, \Sigma_{12}^b, \Sigma_{21}^b, \Sigma_{22}^b \\ \hline \end{array}$$

EM Algorithm

- Tahap Expectation (E):
 - Hitung ekspektasi kelas setiap data

Maximization

- Tahap ~~Expectation (M)~~:
 - Hitung nilai parameter dengan ekspektasi diatas

$$\hat{\gamma}_{ik} = \frac{\hat{\phi}_k \mathcal{N}(x_i | \hat{\mu}_k, \hat{\sigma}_k)}{\sum_{j=1}^K \hat{\phi}_j \mathcal{N}(x_i | \hat{\mu}_j, \hat{\sigma}_j)}$$
$$\hat{\phi}_k = \sum_{i=1}^N \frac{\hat{\gamma}_{ik}}{N}$$
$$\hat{\mu}_k = \frac{\sum_{i=1}^N \hat{\gamma}_{ik} x_i}{\sum_{i=1}^N \hat{\gamma}_{ik}}$$
$$\hat{\sigma}_k^2 = \frac{\sum_{i=1}^N \hat{\gamma}_{ik} (x_i - \hat{\mu}_k)^2}{\sum_{i=1}^N \hat{\gamma}_{ik}}$$

EM Algorithm

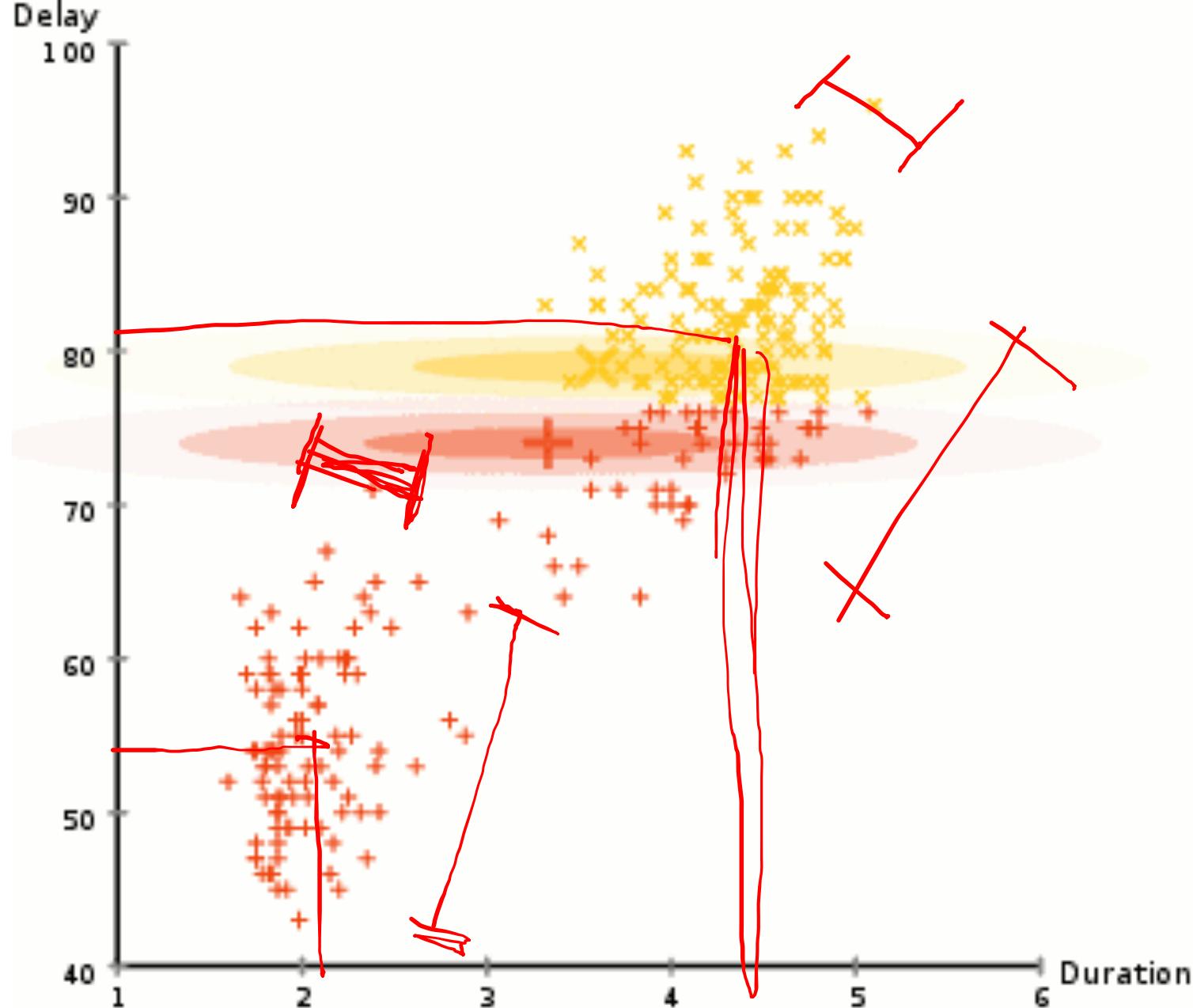
$$D=2, K=2$$

$$\mu_1^r, \mu_2^r, \mu_1^y, \mu_2^y$$

$$\Sigma_{11}^r, \Sigma_{12}^r, \Sigma_{21}^r, \Sigma_{22}^r$$

$$\Sigma_{11}^y, \Sigma_{12}^y, \Sigma_{21}^y, \Sigma_{22}^y$$

$$\pi^r, \pi^y$$



K-means sebagai kasus khusus

- Tahap Expectation (E):
 - Hitung ekspektasi kelas setiap data
- Tahap Expectation (M):
 - Hitung nilai parameter dengan ekspektasi diatas

$$\hat{\gamma}_{ik} = \frac{\hat{\phi}_k \mathcal{N}(x_i | \hat{\mu}_k, \hat{\sigma}_k)}{\sum_{j=1}^K \hat{\phi}_j \mathcal{N}(x_i | \hat{\mu}_j, \hat{\sigma}_j)}$$

$$\hat{\phi}_k = \sum_{i=1}^N \frac{\hat{\gamma}_{ik}}{N}$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^N \hat{\gamma}_{ik} x_i}{\sum_{i=1}^N \hat{\gamma}_{ik}}$$

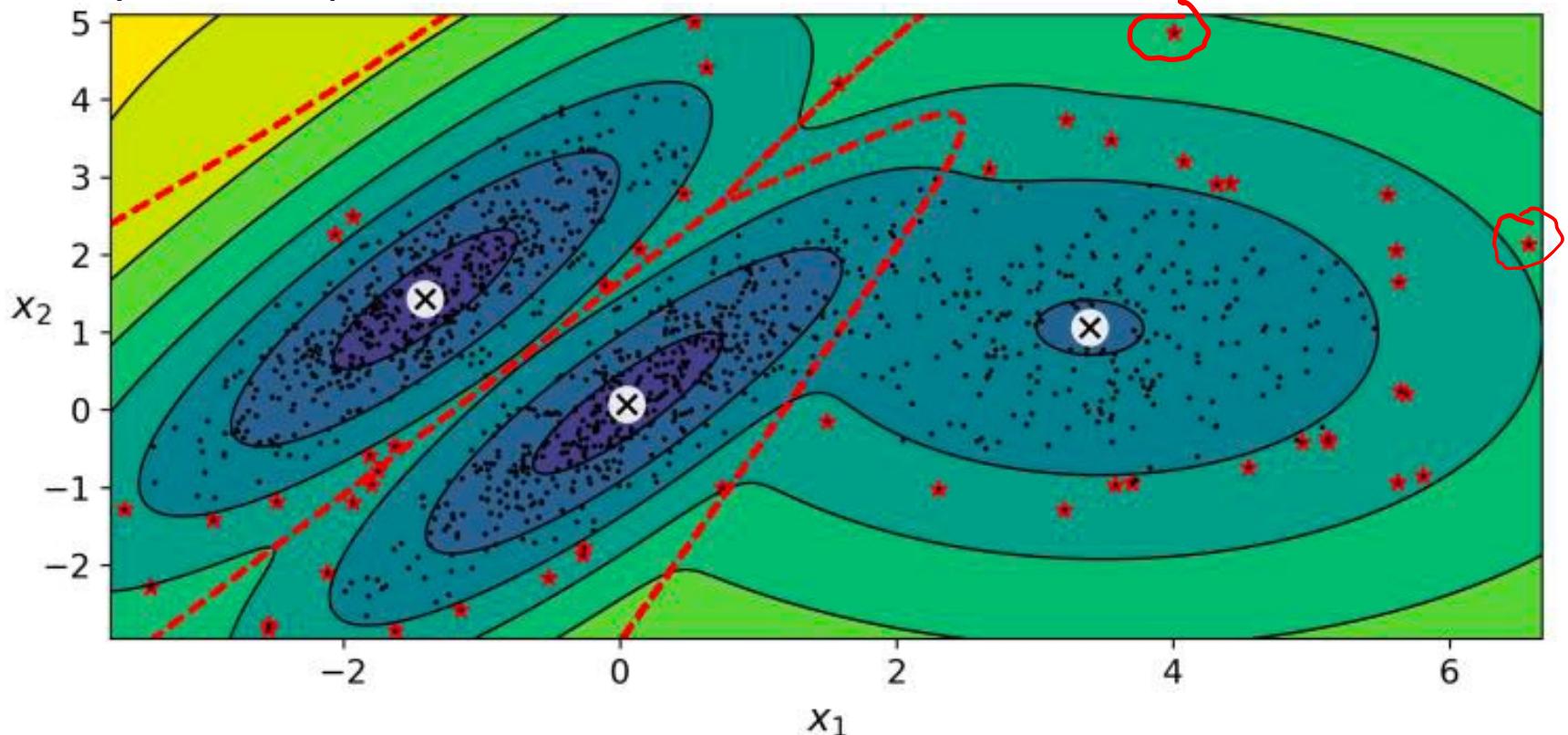
Centroid

$$\hat{\sigma}_k^2 = \frac{\sum_{i=1}^N \hat{\gamma}_{ik} (x_i - \hat{\mu}_k)^2}{\sum_{i=1}^N \hat{\gamma}_{ik}}$$

SSE

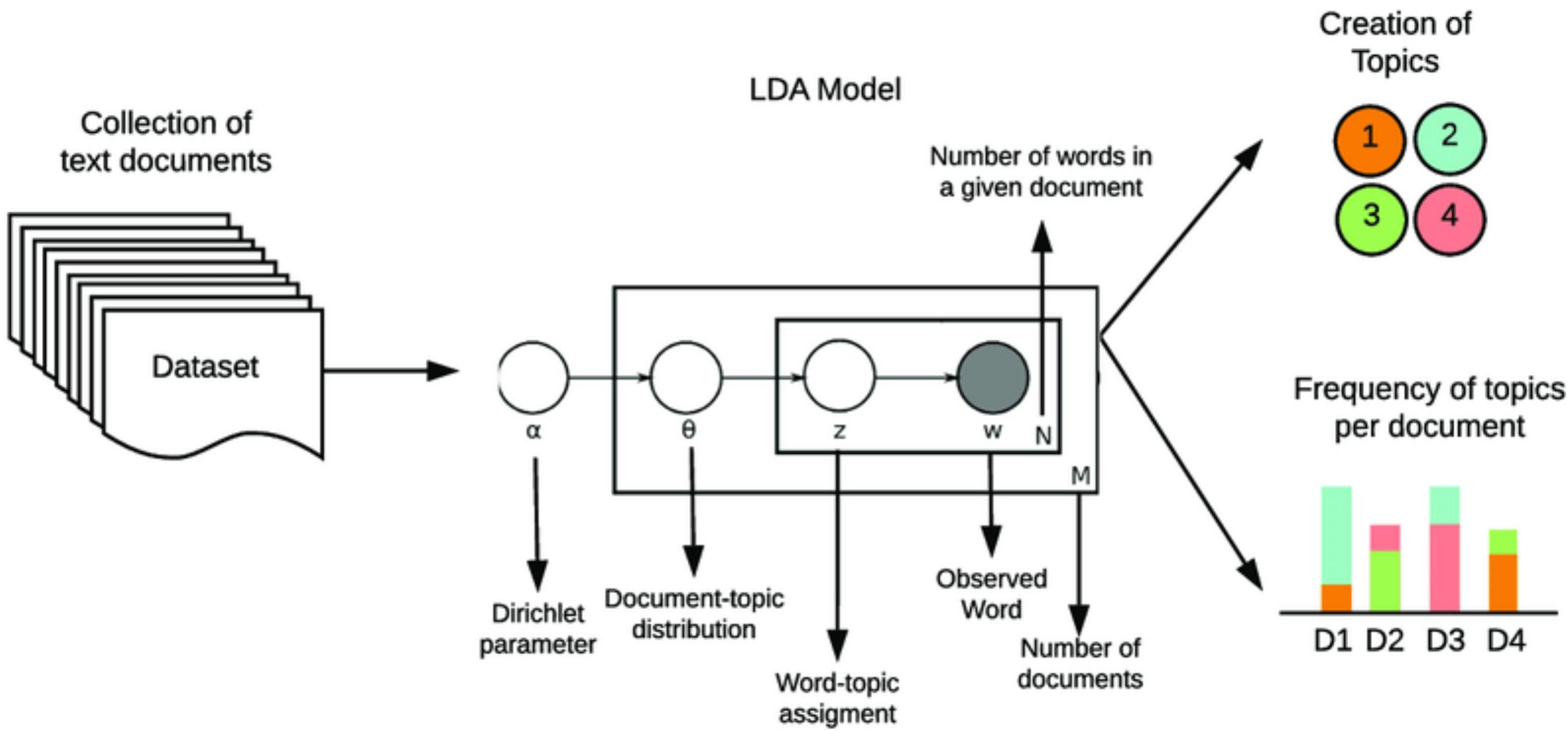
GMM untuk deteksi anomali

- Anomaly: data di area yang tidak padat
- Tentukan density threshold, misal: area kepadatan $< 4\% \rightarrow$ outliers
- Jika terlalu banyak false positive \rightarrow turunkan threshold



Latent Dirichlet Allocation

- Latent Dirichlet Allocation (LDA) adalah algoritma machine learning dari kelompok probabilistic model untuk mempelajari topik tersembunyi dari sebuah teks.
- Latent Dirichlet Allocation berbeda dengan linear discriminant analysis
- Latent Dirichlet Allocation mereduksi bag-of-words menjadi 2 matriks yang bersifat probabilistic yaitu: document-topic matrix dan topic-word matrix
- Matriks probabilistik berarti Jumlah total dari setiap nilai elemen adalah satu
- https://en.wikipedia.org/wiki/Latent_Dirichlet_allocation



Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

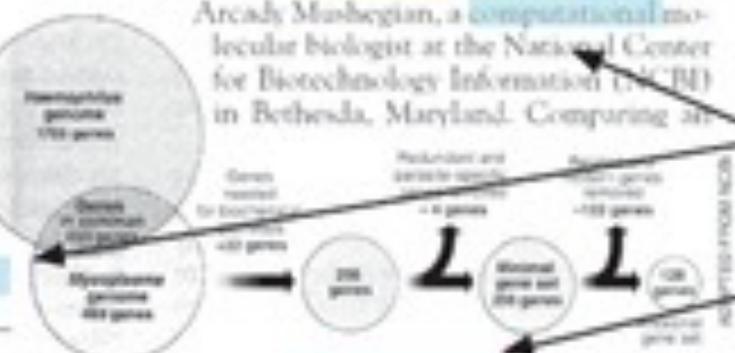
Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an **organism** need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using **computational analyses** to compare known genomes, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a single parasite and estimated that for this organism, 500 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

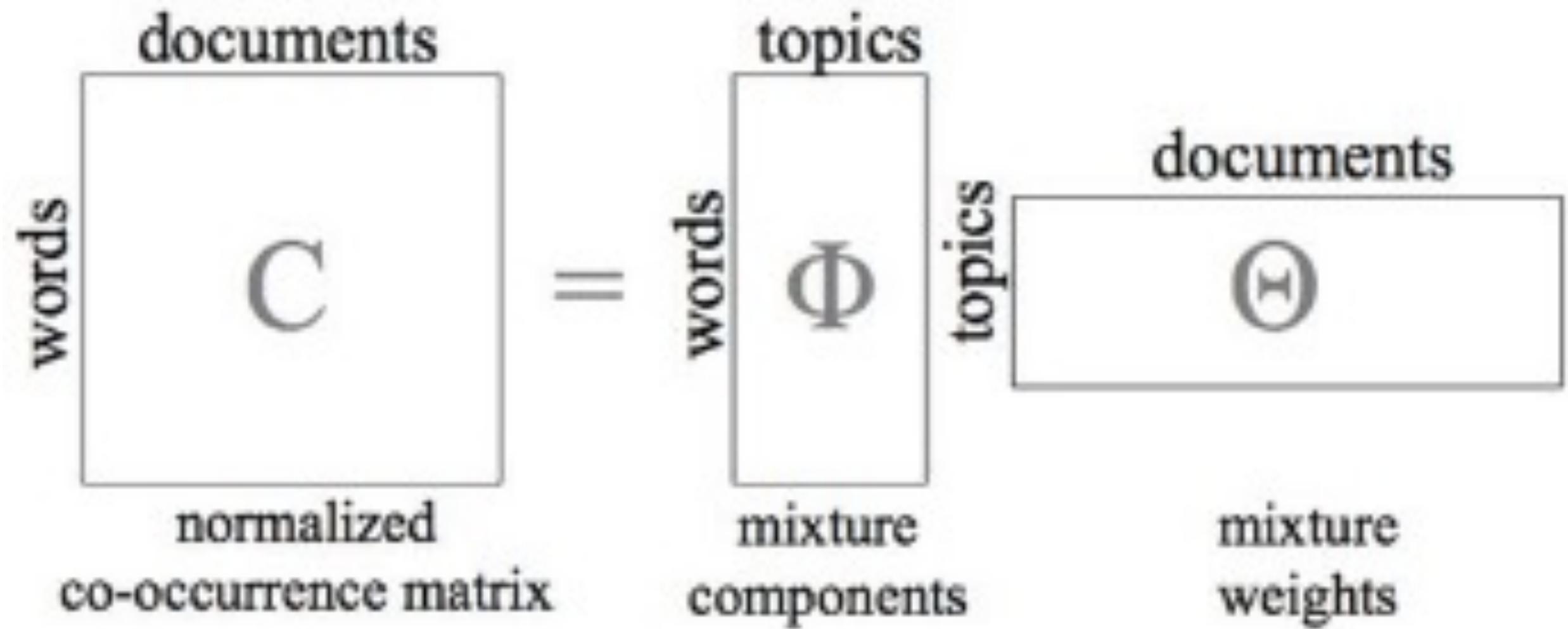
Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Sir Andersson of Umeå University in Sweden, who derived at the 500 number. But coming up with a consensus answer may be more than just a **matter of numbers**. Since, particularly, more and more **genomes** are **continually sequenced** and **sequenced**, "it may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

Topic proportions and assignments



Topic 1	Topic 2	Topic 3	Topic 4
“Technology”	“Entertainment”	“Health”	“U.S. Politics”
company	film	drug	republican
mobile	show	drugs	house
technology	music	cancer	senate
facebook	year	fda	president
google	television	patients	state
apple	singer	reuters	republicans
online	years	disease	political
industry	movie	treatment	campaign
video	band	virus	party
business	actor	health	democratic

Tuhan Memberkati

