

بسمه تعالی

# فصل پنجم داده کاوی

خوشه بندی  
Clustering

مدرس

فاطمه دارائی

[f\\_daraei@semnan.ac.ir](mailto:f_daraei@semnan.ac.ir)

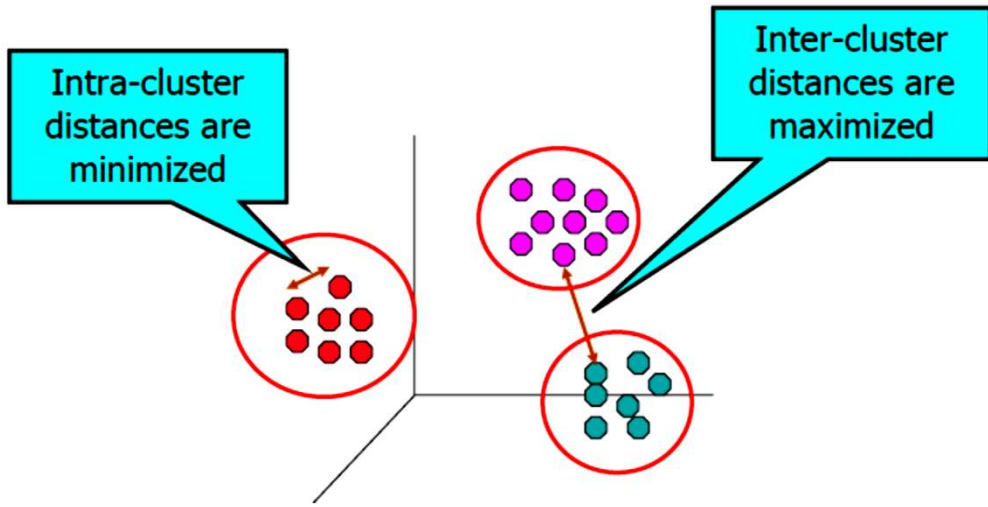
<https://fdaraei.profile.semnan.ac.ir>



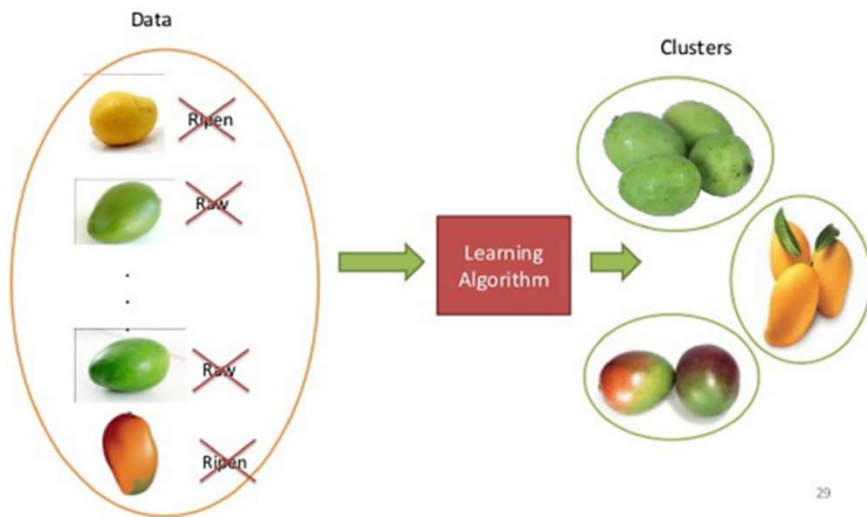


## تحلیل خوشه ای چیست؟

- خوشه‌بندی: مجموعه‌ای از اشیاء داده
- شباهت (یا ارتباط) بین اشیاء درون یک گروه: شباهت **درون کلاسی** بالا
- عدم شباهت (یا عدم ارتباط) با اشیاء گروه‌های دیگر: شباهت **بین کلاسی** پایین
- چه اتفاقی می‌افتد اگر هر یک از این دو شرط را کاهش دهیم؟
- تحلیل خوشه‌ای (یا خوشه‌بندی، تقسیم‌بندی داده‌ها و ...)
- یافتن شباهت‌ها بین داده‌ها براساس ویژگی‌های موجود در داده‌ها
- گروه‌بندی اشیاء مشابه داده به خوشه‌ها



## یادگیری با نظارت Supervised learning



- یادگیری بدون نظارت:
- کلاس‌های از پیش تعریف‌شده وجود ندارد (یعنی یادگیری از طریق مشاهده به جای یادگیری از طریق مثال: یادگیری نظارت‌شده).

### • کاربردهای معمول:

1. به‌عنوان یک ابزار مستقل برای درک توزیع داده‌ها
2. به‌عنوان یک مرحله پیش‌پردازش برای الگوریتم‌های دیگر

<https://www.slideshare.net/sachinnagargoje1/introduction-to-machinelearningatsapthgiricollegebangalore>

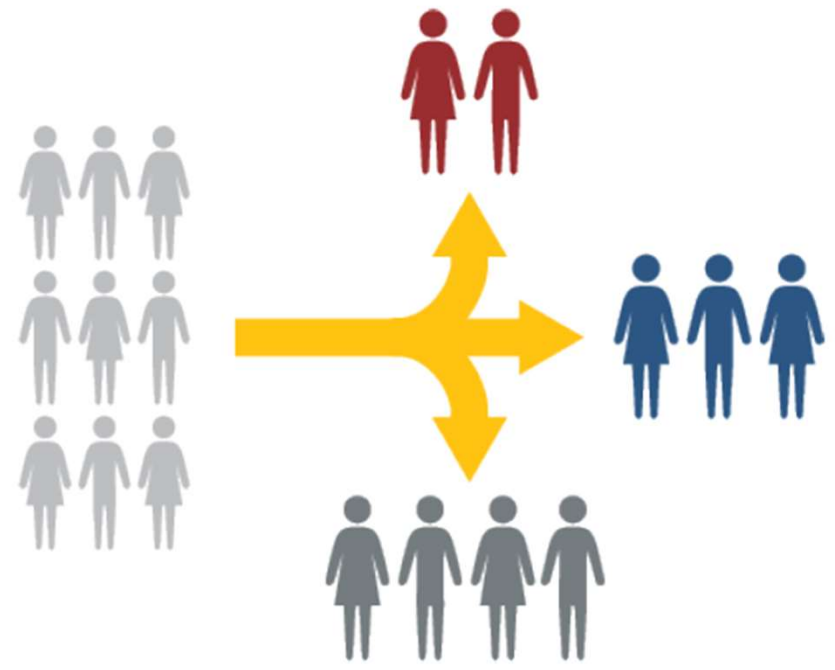


دانشگاه سمنان

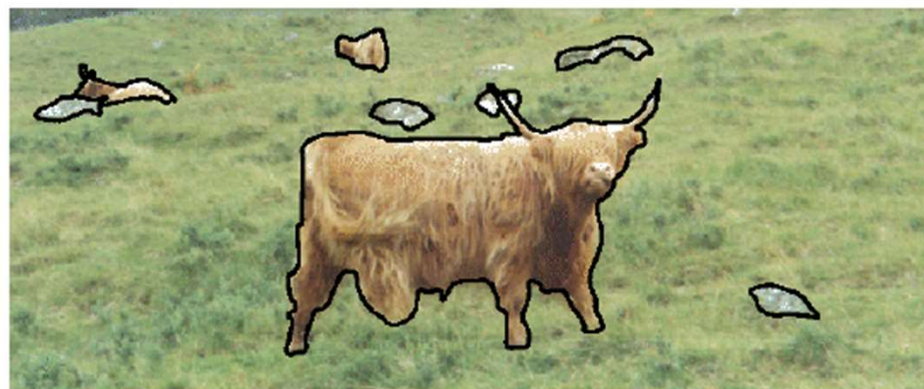
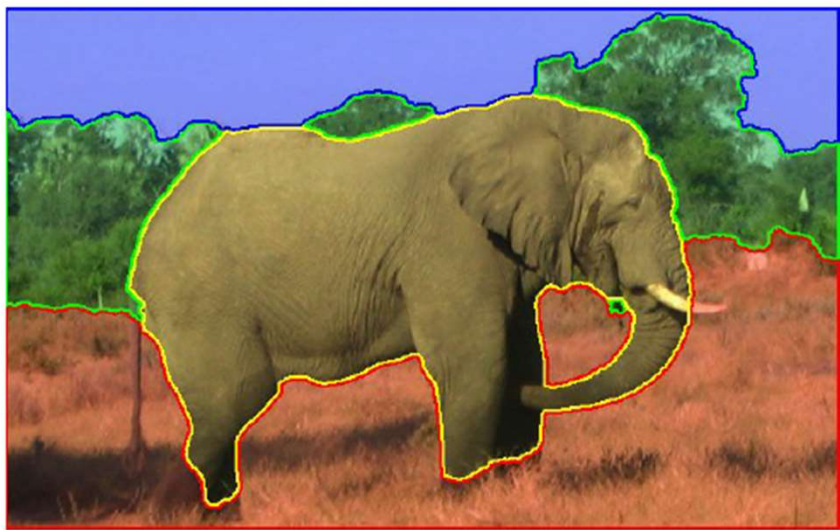
دانشگاه سمنان  
Semnan University

پودیس فرزانتگان

مثال: تقسیم بندی مشتری



مثال: تقسیم بندی تصویر



## مثال: شناسایی داده های پرت

- نقاط پرت **Outliers** چیست؟
- مجموعه‌ای از اشیاء که به طور قابل توجهی با بقیه داده‌ها متفاوت هستند.
- کاربردهای تشخیص نقاط پرت:
- **تشخیص تقلب**، مانند تشخیص تقلب در کارت‌های اعتباری.
- به عنوان یک ابزار **پیش‌پردازش** نیز مفید است.
- یکی از روش‌های تشخیص نقاط پرت با استفاده از خوشه‌بندی:
- اشیائی که به هیچ خوشه‌ای تعلق ندارند.
- اشیائی که فاصله زیادی از سایر اشیاء در همان خوشه دارند.
- خوشه‌هایی با تعداد اعضای بسیار کم



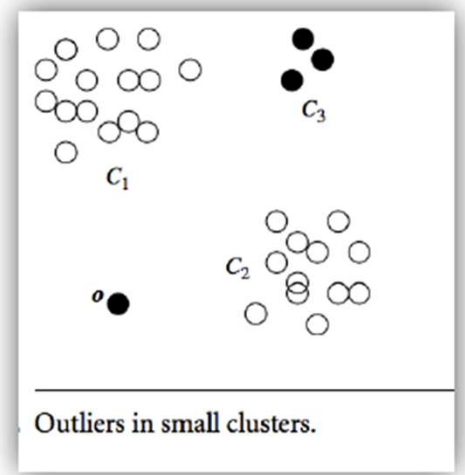
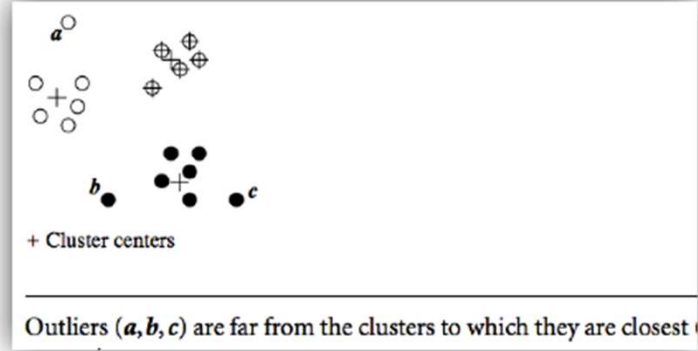
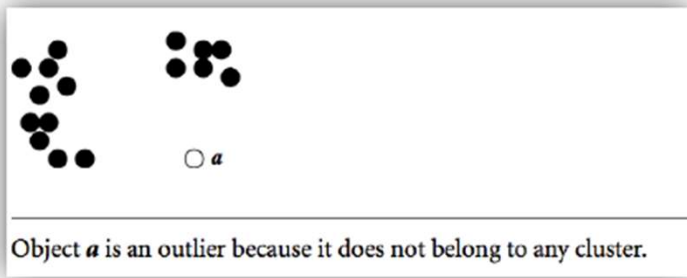
دانشگاه سمنان

دانشگاه سمنان

Semnan University

پودیس فرزانگان

# مثال: شناسایی داده های پرت





دانشگاه سمنان

دانشگاه سمنان

Semnan University

پردیس فرزندانگان

## الگوریتم های Partitioning

نسخه‌های ساده و بنیادی از تحلیل خوشه‌بندی هستند.

اشیای یک مجموعه را به چند گروه یا خوشه‌ی مجزا سازمان‌دهی می‌کند.

خوشه‌ها به گونه‌ای تشکیل می‌شوند که یک معیار هدف برای تقسیم‌بندی بهینه شود.

معیار تقسیم‌بندی: تغییرات درون خوشه‌ای Within cluster variation:

تقسیم یک پایگاه داده  $D$  شامل  $N$  شیء به مجموعه‌ای از  $k$  خوشه، به طوری که مجموع فواصل به توان دو کمینه شود (که در آن  $c_i$  مرکز یا نقطه میانی خوشه  $C_i$  است)

$$E = \sum_{i=1}^k \sum_{p \in C_i} (p - c_i)^2$$





دانشگاه سمنان

دانشگاه سمنان

Semnan University

پردیس فرزانهگان

با توجه به  $k$ ، یک تقسیم‌بندی از  $k$  خوشه پیدا کنید که معیار تقسیم‌بندی انتخاب‌شده را بهینه کند:  
بهینه‌ی کلی Global: تمام تقسیم‌بندی‌ها را به صورت کامل بررسی کنید.

روش‌های اکتشافی Heuristic:

الگوریتم‌های K-means و K-medoids

K-means: خوشه توسط مرکز خوشه نمایش داده می‌شود.

K-medoids تقسیم‌بندی حول میانه‌ها: هر خوشه توسط یکی از اشیای موجود در خوشه نمایش داده می‌شود.

## K-means الگوریتم

**Algorithm:  $k$ -means.** The  $k$ -means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

**Input:**

- $k$ : the number of clusters,
- $D$ : a data set containing  $n$  objects.

**Output:** A set of  $k$  clusters.

**Method:**

- (1) arbitrarily choose  $k$  objects from  $D$  as the initial cluster centers;
- (2) **repeat**
- (3) (re)assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;
- (4) update the cluster means, that is, calculate the mean value of the objects for each cluster;
- (5) **until** no change;



دانشگاه سمنان

دانشگاه سمنان  
Semnan University  
پودیس فرزانگان

# K-means مثال الگوریتم

TABLE 10.1 Data points for  $k$ -means example

<i>A</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>E</i>	<i>f</i>	<i>g</i>	<i>h</i>
(1,3)	(3,3)	(4,3)	(5,3)	(1,2)	(4,2)	(1,1)	(2,1)

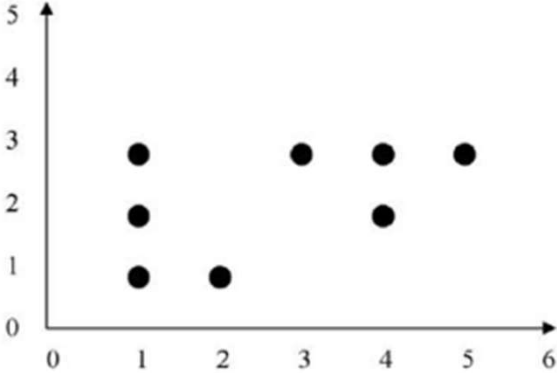


Figure 10.4 How will  $k$ -means partition these data into  $k = 2$  clusters?

# حل مثال الگوریتم K-means

**TABLE 10.1 Data points for *k*-means example**

<i>A</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>E</i>	<i>f</i>	<i>g</i>	<i>h</i>
(1,3)	(3,3)	(4,3)	(5,3)	(1,2)	(4,2)	(1,1)	(2,1)

$m_1=(1,1)$   
 $m_2=(2,1)$

**TABLE 10.2 Finding the nearest cluster center for each record (first pass)**

Point	Distance from $m_1$	Distance from $m_2$	Cluster Membership
<i>a</i>	2.00	2.24	$C_1$
<i>b</i>	2.83	2.24	$C_2$
<i>c</i>	3.61	2.83	$C_2$
<i>d</i>	4.47	3.61	$C_2$
<i>e</i>	1.00	1.41	$C_1$
<i>f</i>	3.16	2.24	$C_2$
<i>g</i>	0.00	1.00	$C_1$
<i>h</i>	1.00	0.00	$C_2$

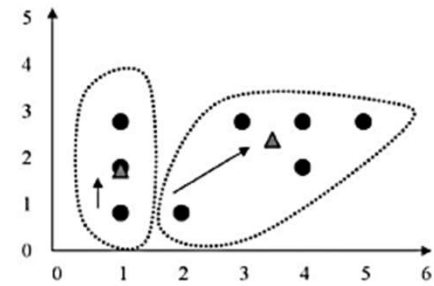


Figure 10.5 Clusters and centroids  $\Delta$  after first pass through *k*-means algorithm.

## ادامه حل مثال الگوریتم K-means

**TABLE 10.1** Data points for *k*-means example

<i>A</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>E</i>	<i>f</i>	<i>g</i>	<i>h</i>
(1,3)	(3,3)	(4,3)	(5,3)	(1,2)	(4,2)	(1,1)	(2,1)

$$m_1 = (1, 2)$$

$$m_2 = (3.6, 2.4)$$

**TABLE 10.3** Finding the nearest cluster center for each record (second pass)

Point	Distance from $m_1$	Distance from $m_2$	Cluster Membership
<i>a</i>	1.00	2.67	$C_1$
<i>b</i>	2.24	0.85	$C_2$
<i>c</i>	3.16	0.72	$C_2$
<i>d</i>	4.12	1.52	$C_2$
<i>e</i>	0.00	2.63	$C_1$
<i>f</i>	3.00	0.57	$C_2$
<i>g</i>	1.00	2.95	$C_1$
<i>h</i>	1.41	2.13	$C_1$

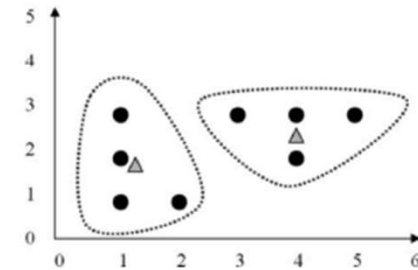


Figure 10.6 Clusters and centroids  $\Delta$  after second pass through *k*-means algorithm.

## K-means ادامه حل مثال الگوریتم

**TABLE 10.1 Data points for  $k$ -means example**

<i>A</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>E</i>	<i>f</i>	<i>g</i>	<i>h</i>
(1,3)	(3,3)	(4,3)	(5,3)	(1,2)	(4,2)	(1,1)	(2,1)

$$m_1 = (1.25, 1.75)$$

$$m_2 = (4, 2.75)$$

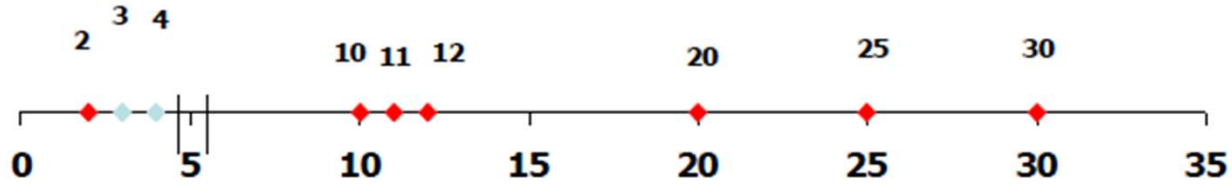
**TABLE 10.4 Finding the nearest cluster center for each record (third pass)**

Point	Distance from $m_1$	Distance from $m_2$	Cluster Membership
<i>a</i>	1.27	3.01	$C_1$
<i>b</i>	2.15	1.03	$C_2$
<i>c</i>	3.02	0.25	$C_2$
<i>d</i>	3.95	1.03	$C_2$
<i>e</i>	0.35	3.09	$C_1$
<i>f</i>	2.76	0.75	$C_2$
<i>g</i>	0.79	3.47	$C_1$
<i>h</i>	1.06	2.66	$C_1$

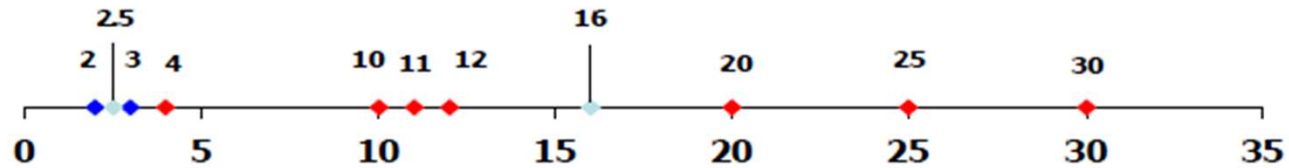
No change in clustering

## مثال الگوریتم K-means

- $\{2,4,10,12,3,20,30,11,25\}$ ,  $k=2$
- $m_1=3, m_2=4$

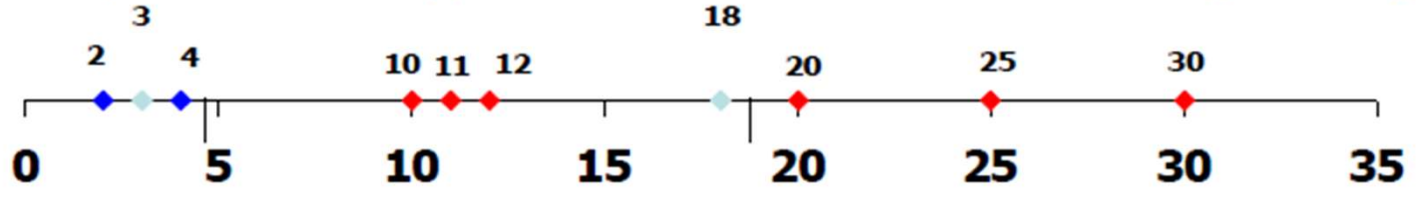


- $K_1=\{2,3\}$ ,  $K_2=\{4,10,12,20,30,11,25\}$ ,  $m_1=2.5, m_2=16$

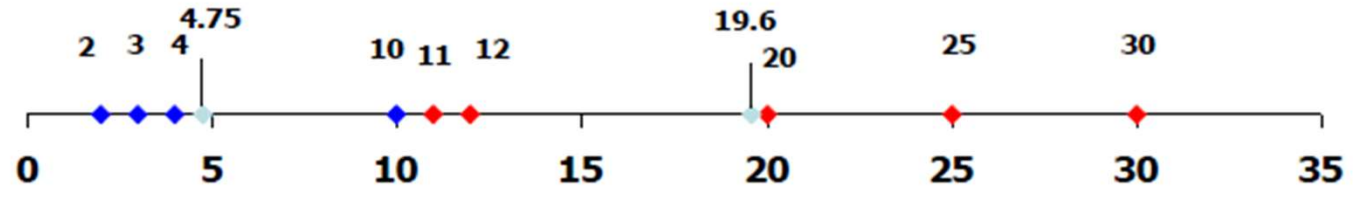




- $K_1 = \{2, 3, 4\}, K_2 = \{10, 12, 20, 30, 11, 25\}, m_1 = 3, m_2 = 18$

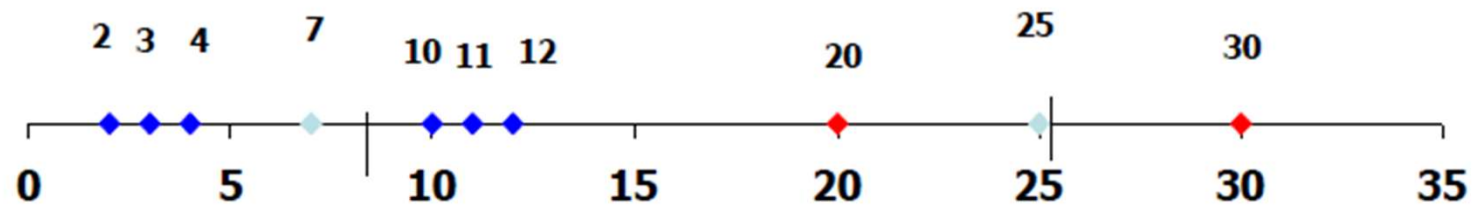


- $K_1 = \{2, 3, 4, 10\}, K_2 = \{12, 20, 30, 11, 25\}, m_1 = 4.75, m_2 = 19.6$





- $K_1 = \{2, 3, 4, 10, 11, 12\}, K_2 = \{20, 30, 25\}, m_1 = 7, m_2 = 25$



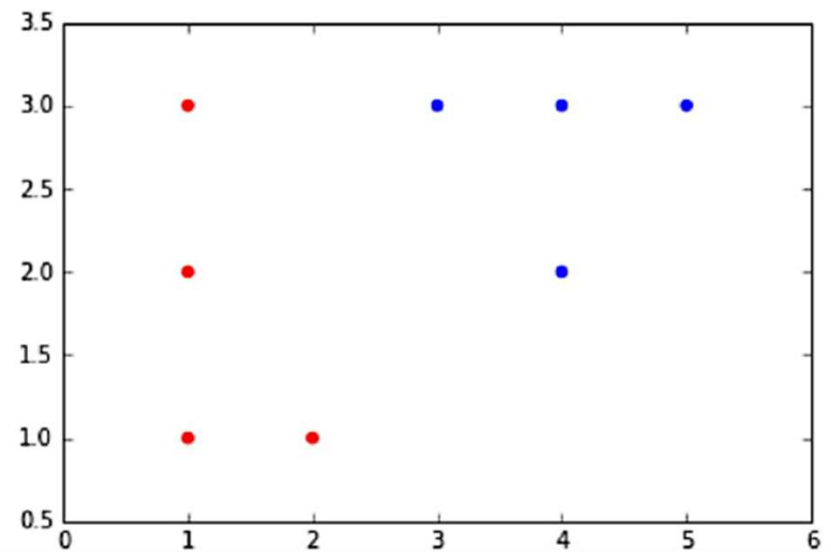
```
from sklearn.cluster import KMeans
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

X = np.array([[1, 3], [3, 3], [4, 3], [5, 3],
              [1, 2], [4, 2], [1, 1], [2, 1]])

kmeans = KMeans(n_clusters=2).fit(X)
print(kmeans.labels_)

colors = np.array(['r', 'b'])
plt.scatter(X[:,0], X[:,1], color=colors[kmeans.labels_])

[0 1 1 1 0 1 0 0]
```





## K-means

مزایا:

کارآمد و مناسب برای داده های بزرگ: از مرتبه  $O(tkn)$ ، که در آن  $n$  تعداد اشیا،  $k$  تعداد خوشه ها، و  $t$  تعداد تکرارها است. معمولاً  $k, t \ll n$

معایب:

اغلب در یک نقطه بهینه محلی متوقف می شود.

فقط برای اشیایی در فضای پیوسته  $n$ -بعدی قابل استفاده است.

برای داده های دسته بندی شده باید از روش  $k$ -modes استفاده شود.

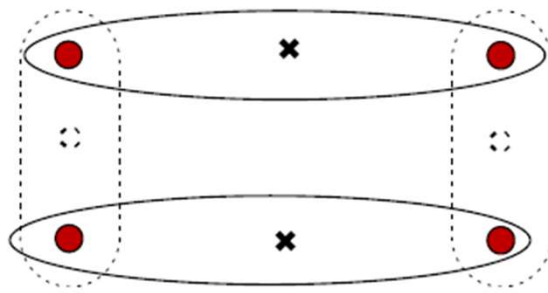
در مقایسه، میانه ها را می توان برای طیف گسترده ای از داده ها اعمال کرد.

نیاز به مشخص کردن  $k$ ، تعداد خوشه ها، از پیش دارد.

حساس به داده های نویزی و نقاط پرت: یک راه حل  $k$ -medoids است.

حساس به خوشه های اولیه: یک راه حل  $k$ -means++ است.

برای کشف خوشه های با اشکال غیر محدب مناسب نیست.



k-means حساس به خوشه‌های اولیه است.

با انتخاب دقیق مراکز اولیه خوشه‌ها، ممکن است بتوانیم نه تنها سرعت همگرایی الگوریتم را افزایش دهیم، بلکه کیفیت خوشه‌بندی نهایی را نیز تضمین کنیم. الگوریتم `k-means++` یک نسخه از `k-means` است که مراکز اولیه را به شرح زیر انتخاب می‌کند:

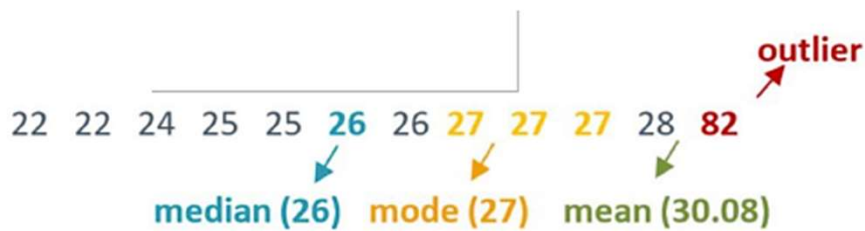
ابتدا، یک مرکز به طور تصادفی از اشیای موجود در مجموعه داده انتخاب می‌شود. مراکز دیگر به طور تصادفی و با احتمال متناسب با  $\text{dist}(p)^2$  انتخاب می‌شوند، جایی که  $\text{dist}(p)$  فاصله نقطه  $p$  از نزدیک‌ترین مرکزی است که قبلاً انتخاب شده است.

در `sklearn.cluster.KMeans`، روش اولیه پیش‌فرض، `k-means++` است.

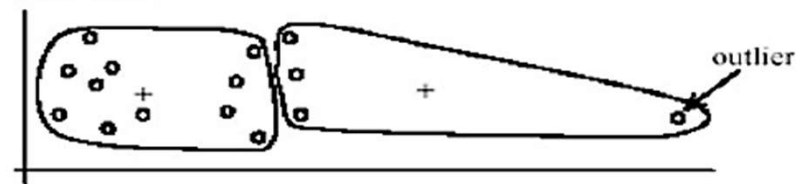
# k-medoids

الگوریتم  $k$ -means به داده‌های پرت حساس است!

$k$ -medoids: به جای استفاده از میانگین اشیاء در یک خوشه به عنوان نقطه مرجع، می‌توان از میانه‌ها استفاده کرد، که اشیاء با موقعیت مرکزی‌ترین در یک خوشه هستند.



<https://www.cese.nsw.gov.au/effective-practices/unit-4-outliers>



(A): Undesirable clusters



(B): Ideal clusters

<https://www.slideshare.net/anilyadav5055/15857-cse422-unsupervisedlearning>

خوشه بندی K-Medoids: یافتن اشیاء نماینده (Medoids) در خوشه‌ها

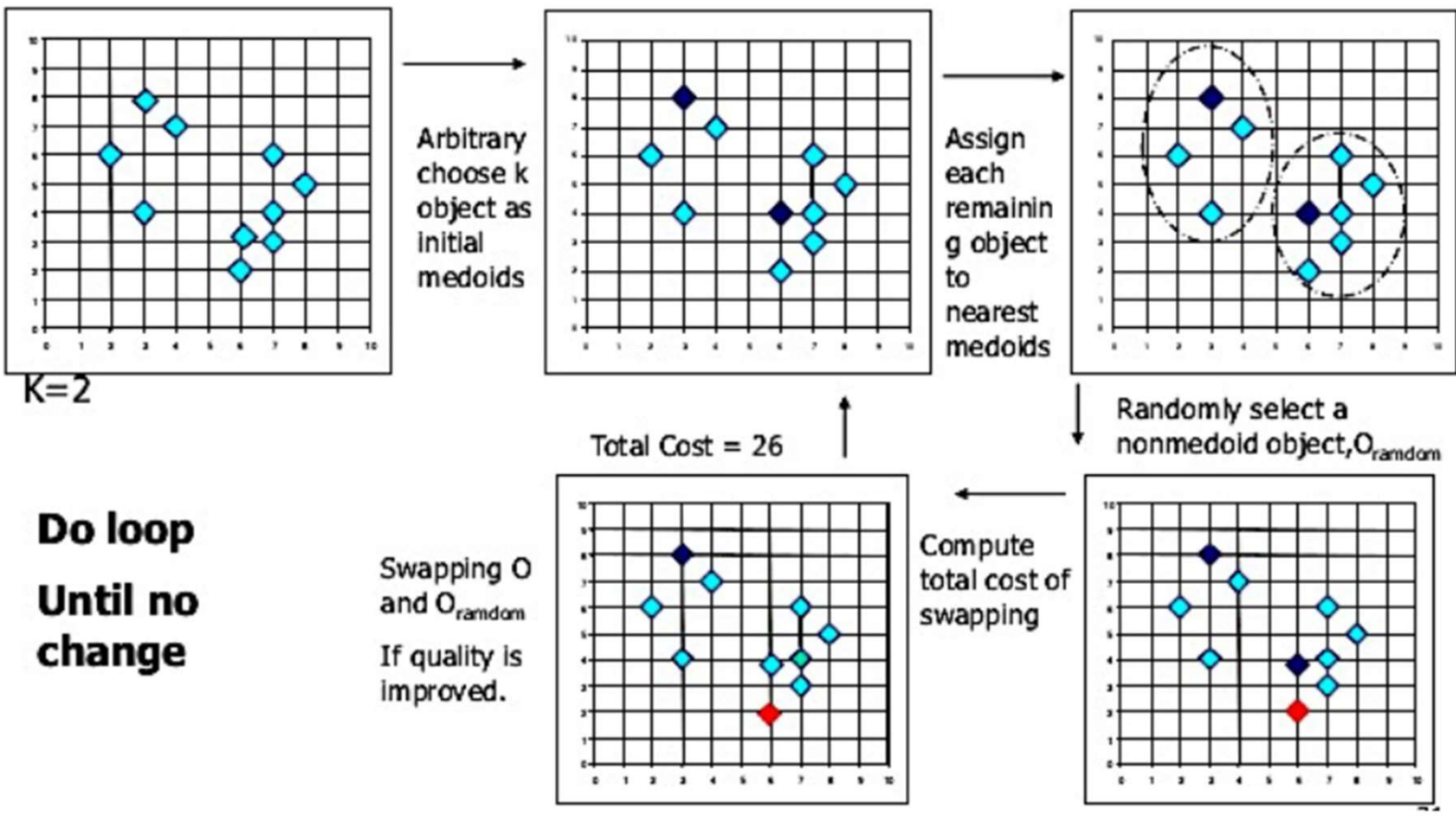
PAM (Partitioning Around Medoids, Kaufmann & Rousseeuw 1987):

از مجموعه اولیه‌ای از میانه‌ها شروع کرده و به صورت تکراری یکی از Medoids را با یکی از non-Medoids جایگزین می‌کند اگر که کل مسافت خوشه‌بندی بهبود یابد.

PAM برای مجموعه‌های داده کوچک مؤثر است، اما برای مجموعه‌های داده بزرگ مقیاس‌پذیری ندارد (به دلیل پیچیدگی محاسباتی)

## الگوریتم k-medoids

- شیء به عنوان medoid های اولیه به صورت دلخواه اختیار کن.
- تکرار کن تا اینکه هیچ تغییری رخ ندهد.
- هر کدام از اشیاء باقیمانده را به خوشه‌ای با نزدیکترین medoid تخصیص بده
- بطور تصادفی یک شی غیر medoid را انتخاب کن .،
- هزینه نهایی S را از عوض کردن ( medoid آن خوشه) و محاسبه کن
- اگر  $s < 0$  آنگاه جای عناصر را عوض کن تا مجموعه K تا medoid جدید شکل بگیرد.



**Do loop  
Until no  
change**



## مثال k-medoids

- **1, 2, 6, 7, 8, 10, 15, 17, 20 – break into 3 clusters**
  - Cluster = 6 – 1, 2
  - Cluster = 7
  - Cluster = 8 – 10, 15, 17, 20
- **Random non-medoid – 15 replace 7 (total cost=-13)**
  - Cluster = 6 – 1 (cost 0), 2 (cost 0), 7(1-0=1)
  - Cluster = 8 – 10 (cost 0)
  - New Cluster = 15 – 17 (cost 2-9=-7), 20 (cost 5-12=-7)
- **Replace medoid 7 with new medoid (15) and reassign**
  - Cluster = 6 – 1, 2, 7
  - Cluster = 8 – 10
  - Cluster = 15 – 17, 20

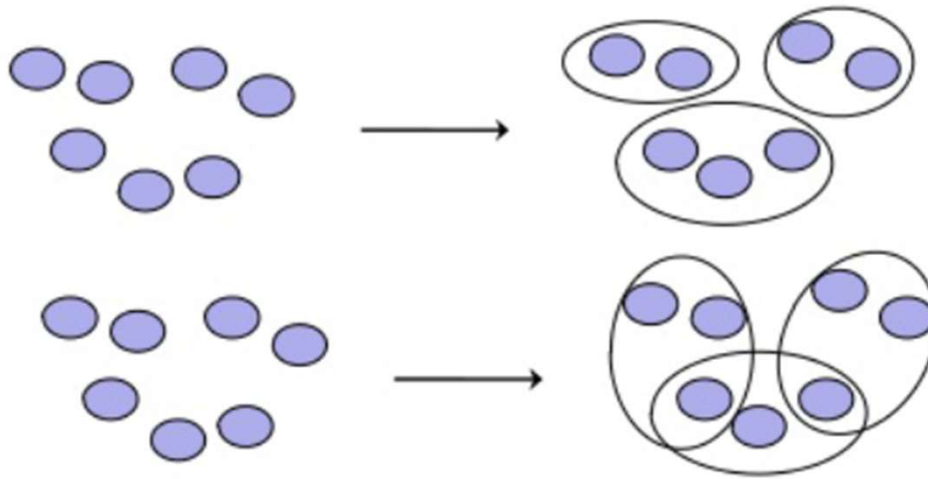
- **Random non-medoid – 1 replaces 6 (total cost= -1)**
  - Cluster = 8 – 6 (cost  $2-0=2$ ), 7 (cost  $1-1=0$ ), 10 (cost 0)
  - Cluster = 15 – 17(cost 0), 20(cost 0)
  - New Cluster = 1 – 2 (cost  $1-4= -3$ )
- **Replace medoid 6 with new medoid (1) and reassign**
  - Cluster = 1 – 2
  - Cluster = 8 – 6, 7, 10
  - Cluster = 15 – 17, 20
- **Random non-medoid – 10 replaces 8 (total cost=2) don't replace**
  - Cluster = 1– 2(cost 0)
  - Cluster = 15 – 17 (cost 0), 20(cost 0)
  - New Cluster = 10 – 6 (cost 0), 7 (cost 0), 8 (cost  $2-0=2$ )

- **Random non-medoid – 17 replaces 15 (total cost=0) don't replace**
  - Cluster = 1 – 2(cost 0)
  - Cluster = 8 – 6 (cost 0), 7 (cost 0), 10 (cost 0)
  - New Cluster = 17 – 15 (cost 2-0=2), 20(cost 3-5=-2)
- **Random non-medoid – 20 replaces 15 (total cost=6) don't replace**
  - Cluster = 1 – 2(cost 0)
  - Cluster = 8 – 6 (cost 0), 7 (cost 0), 10 (cost 0)
  - New Cluster = 20 – 15 (cost 5-0=5), 17(cost 3-2=1)
- **Other possible changes all have high costs**
  - 1 replaces 15, 2 replaces 15, 1 replaces 8, ...
- **No changes, final clusters**
  - Cluster = 1 – 2
  - Cluster = 8 – 6, 7, 10
  - Cluster = 15 – 17, 20

## خوشه بندی فازی یا نرم (Soft (fuzzy) clustering

**خوشه بندی سخت Hard clustering:** هر شیء می تواند دقیقاً به یک خوشه تعلق داشته باشد.  
**خوشه بندی نرم Soft clustering:** اشیاء می توانند با یک درجه عضویت کسری به چندین خوشه تعلق داشته باشند.  
مانند الگوریتم فازی

Fuzzy c-means (FCM)

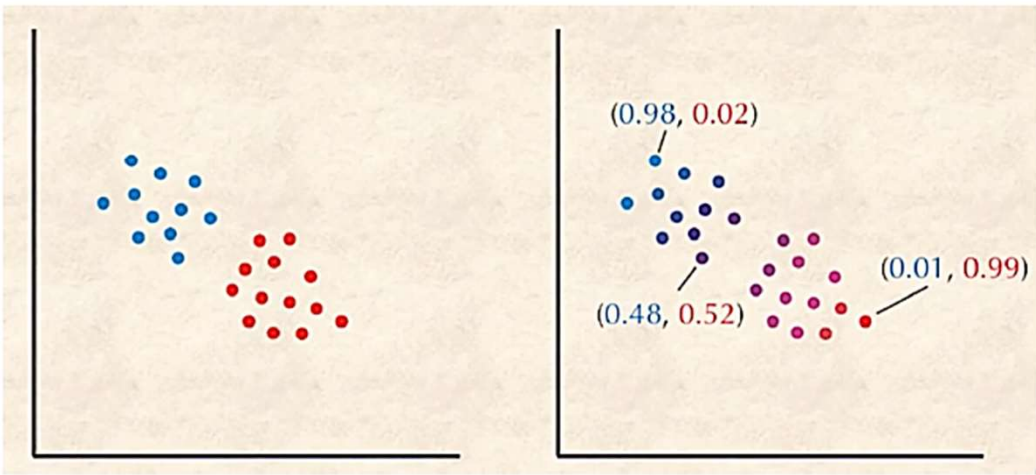
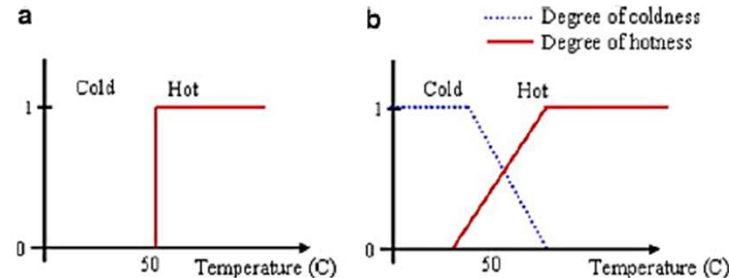


<https://www.slideshare.net/kanimozhiu/text-clustering>

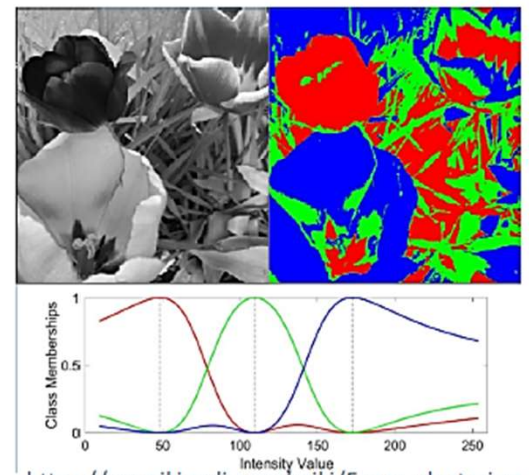
# Soft (fuzzy) clustering



An autobiography of Prof. Zadeh:  
"My life and work—a retrospective view." *Applied and Computational Mathematics* 10.1 (2011): 4-9.



<https://www.youtube.com/watch?v=xtDMHPVDDKk>



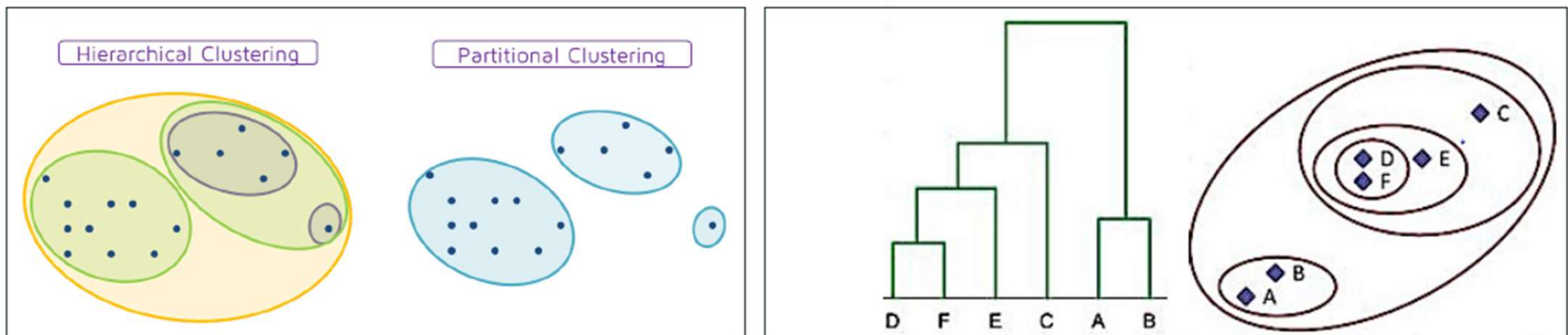
[https://en.wikipedia.org/wiki/Fuzzy\\_clustering](https://en.wikipedia.org/wiki/Fuzzy_clustering)

## خوشه بندی سلسله مراتبی

داده‌ها به چندین سطح از تقسیم‌بندی‌های تو در تو (درخت خوشه‌ها) تجزیه شده، که به نمودار درختی تشکیل شده دندروگرام dendrogram گفته می‌شود.

• نباید هر تعداد خاصی از خوشه را فرض کنیم.

هر تعداد مورد نظر از خوشه را می‌توان با برش dendrogram در سطح مناسب به دست آورد.

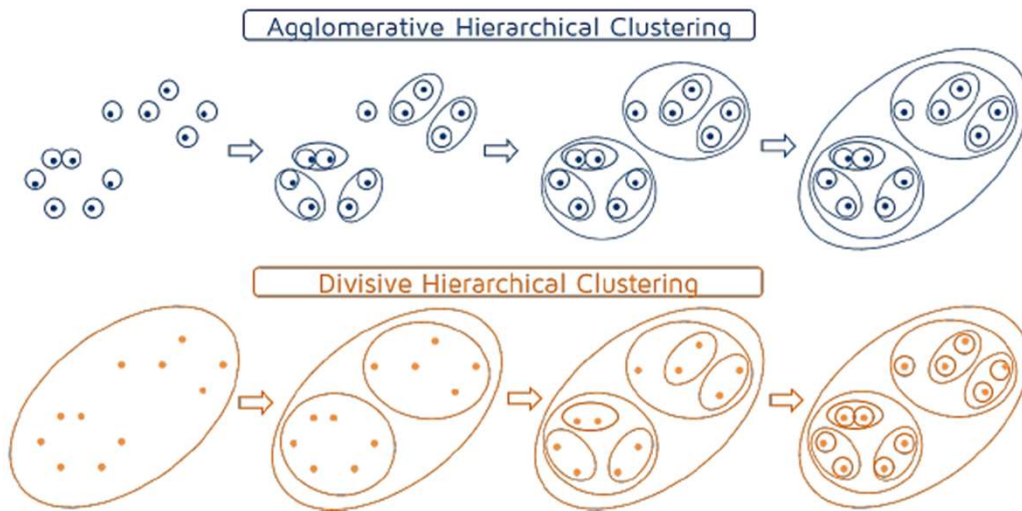


<https://quantdare.com/hierarchical-clustering/>

## روش های خوشه بندی سلسله مراتبی

### تجمیعی Agglomerative (پایین به بالا):

شروع با هر نمونه به عنوان یک خوشه منفرد، می باشد. سپس خوشه ها را با ترکیب به خوشه های بزرگ و بزرگتر تبدیل می کند تا در نهایت، تمام نمونه ها به یک خوشه مشترک تعلق می گیرند.



<https://quantdare.com/hierarchical-clustering/>

### تقسیمی Divisive (بالا به پایین):

شروع با تمام نمونه ها که به یک خوشه مشترک تعلق دارند، می باشد. الگوریتم این خوشه را به خوشه های کوچک و کوچک تر تقسیم می کند تا اینکه هر شی در یک خوشه قرار گیرد. در نهایت، هر نمونه یک خوشه جداگانه تشکیل می دهد.

این روش به دلیل پیچیدگی محاسباتی بالا کمتر مورد استفاده قرار می گیرد.

## فاصله خوشه بندی سلسله مراتبی

الگوریتم‌های سلسله مراتبی سنتی از یک ماتریس شباهت یا فاصله استفاده می‌کنند. روش محاسبه فاصله مهم است.

معیارهای گوناگونی که در روش سلسله مراتبی برای فاصله بین خوشه ها بکار می رود عبارتند از:

**پیوند تک Single Link** : کمترین فاصله بین یک عنصر در یک خوشه و یک عنصر در خوشه دیگر.

**پیوند کامل Complete Link** : بیشترین فاصله بین یک عنصر در یک خوشه و یک عنصر در خوشه دیگر.

**میانگین Average** : میانگین فاصله بین یک عنصر در یک خوشه و یک عنصر در خوشه دیگر.

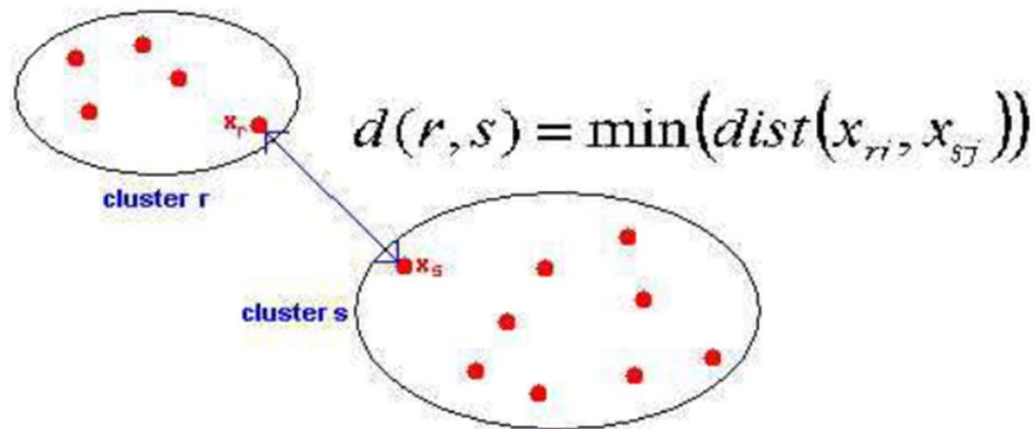
**مرکز ثقل Centroid** : فاصله بین مراکز ثقل دو خوشه.

**مدوید Medoid** : فاصله بین مدویدهای دو خوشه



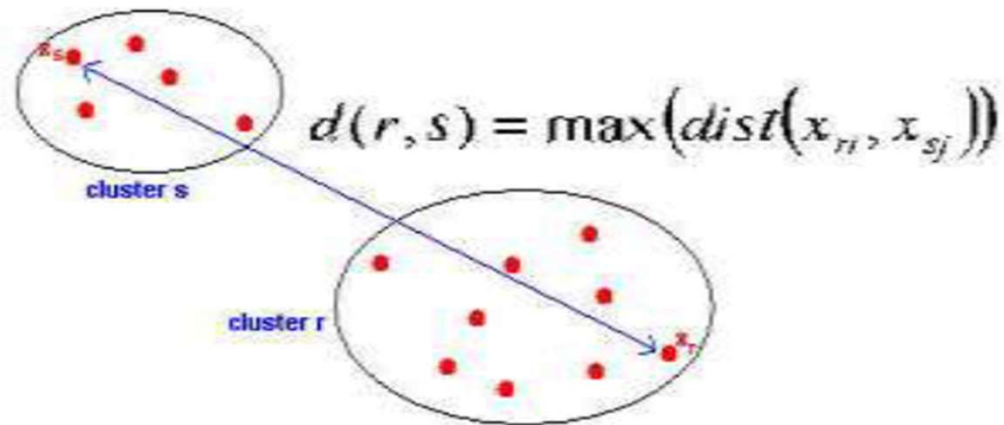
# Single linkage

- فاصله بین دسته‌ها بر حسب حداقل فاصله ممکنه بین عناصر آنها محاسبه می‌شود.
- کلیه فاصله بین زوجهای عناصر دو دسته محاسبه شده و حداقل آنها فاصله بین دو دسته را تعیین می‌کند.



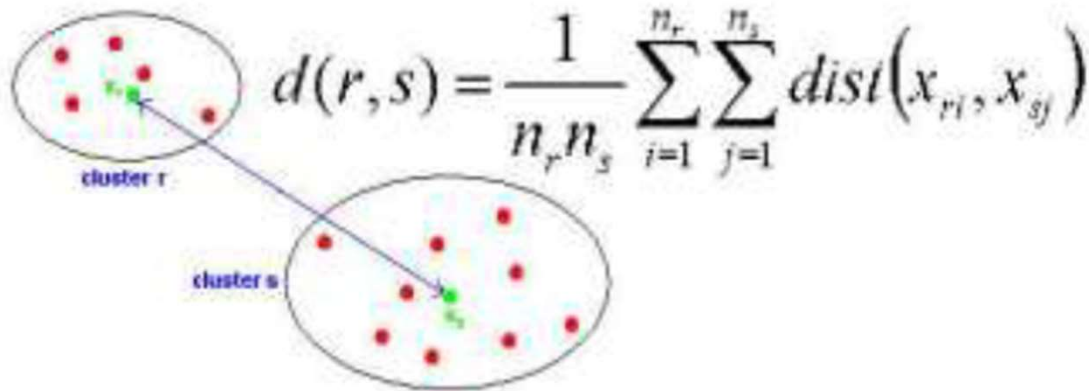
# Complete linkage

- فاصله بین دسته‌ها بر حسب دورترین فاصله ممکنه بین عناصر آنها محاسبه می‌شود.



# Average linkage

- فاصله بین دو دسته مساوی مقادیر متوسط کلیه فاصله‌های ممکنه بین عناصر دو دسته است.


$$d(r, s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} \text{dist}(x_{ri}, x_{sj})$$



دانشگاه سمنان

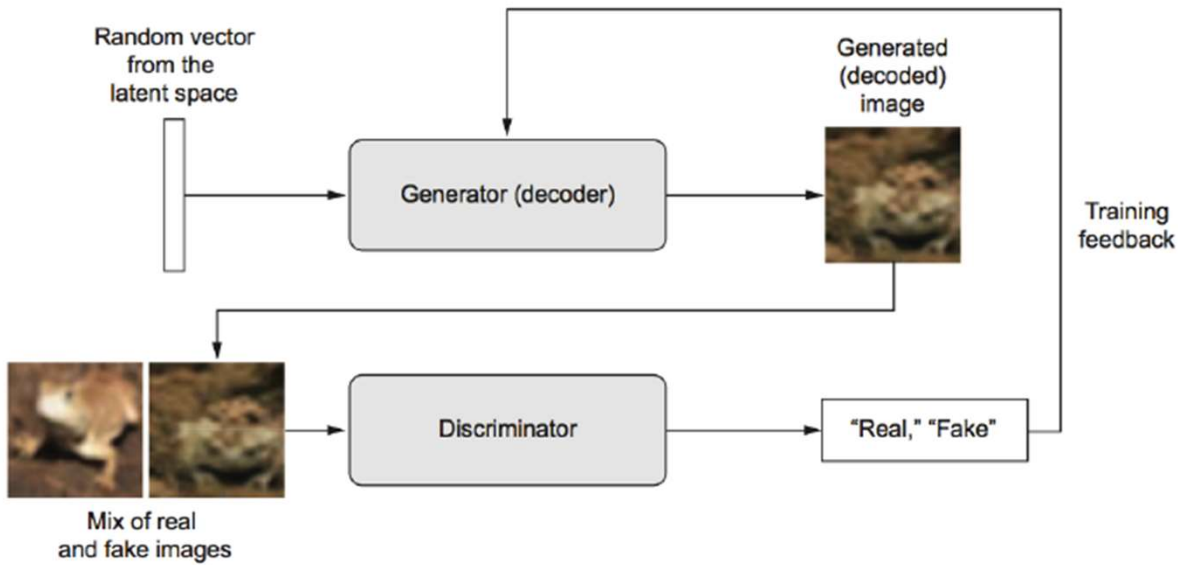
دانشگاه سمنان

Semnan University

پردیس فرزنانگان

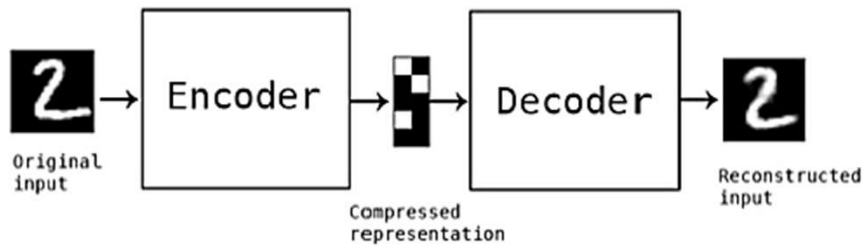
## یادگیری بدون نظارت GAN

- Generative adversarial networks (GANs), introduced in 2014 by Goodfellow et al.

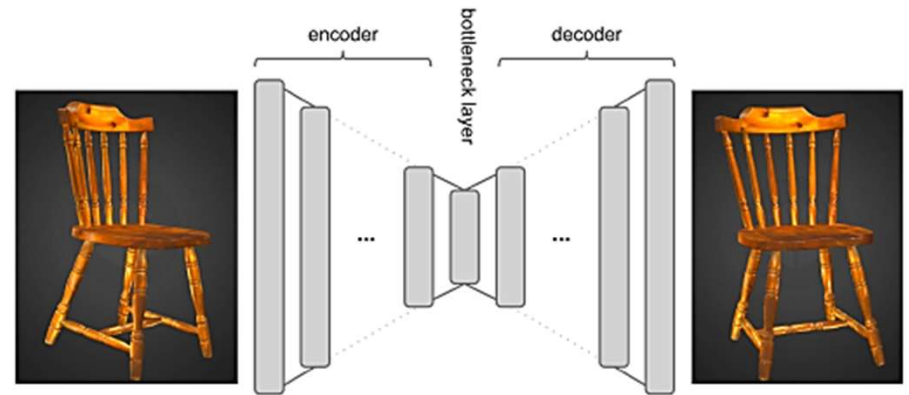


[www.miketyka.com](http://www.miketyka.com)

## یادگیری بدون نظارت Auto encoder



<https://blog.keras.io/building-autoencoders-in-keras.html>



<http://www.inference.vc/stereovision-autoencoder/>

## ارزیابی کیفیت خوشه بندی

معیارهای ارزیابی خوشه بندی:

### ۱. ارزیابی خارجی **Extrinsic**: وجود داده‌های واقعی **Ground Truth**

در این روش، خوشه بندی با استفاده از داده‌های واقعی (**Ground Truth**) مقایسه می‌شود. از معیارهای خاصی برای ارزیابی کیفیت خوشه بندی استفاده می‌شود.

### ۲. ارزیابی داخلی **Intrinsic**: عدم وجود داده‌های واقعی **Ground Truth**

در این روش، بدون استفاده از داده‌های واقعی، کیفیت خوشه بندی بر اساس معیارهایی مانند جدایی خوشه‌ها (**Separability**) و تراکم یا فشردگی خوشه‌ها (**Compactness**) بررسی می‌شود

## ارزیابی خارجی: خلوص Purity

$$\text{purity}(\Omega, C) = \frac{1}{N} \sum_k \max_j |\omega_k \cap c_j|$$

- $\Omega = \{\omega_1, \dots, \omega_k\}$  is the set of clusters and  $C = \{c_1, \dots, c_j\}$  is the set of classes.
- For each cluster  $\omega_k$  : find class  $c_j$  with most members  $n_{kj}$  in  $\omega_k$ .
- Sum all  $n_{kj}$  and divide by total number of points.

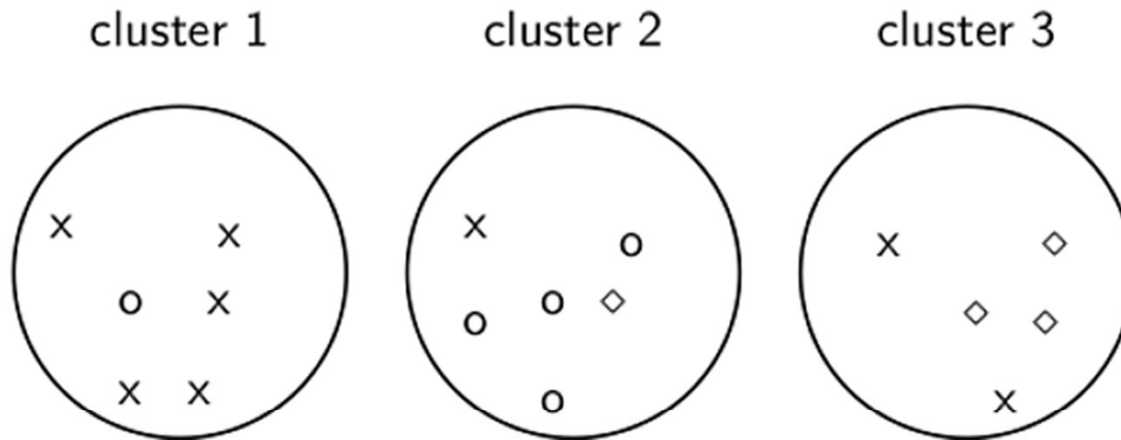


دانشگاه سمنان

دانشگاه سمنان

Semnan University

پردیس فرزندان



$$5 = \max_j |\omega_1 \cap c_j| \text{ (class x, cluster 1);}$$

$$4 = \max_j |\omega_2 \cap c_j| \text{ (class o, cluster 2);}$$

$$3 = \max_j |\omega_3 \cap c_j| \text{ (class } \diamond, \text{ cluster 3).}$$

$$\text{Purity is } (1/17) \times (5 + 4 + 3) \approx 0.71.$$



## ارزیابی خارجی: Rand index

بر اساس جدول تداخل 2x2 برای تمام جفت‌های نمونه‌ها:

در این روش، ارزیابی خوشه‌بندی بر اساس یک جدول تداخل 2x2 Contingency Table انجام می‌شود. این جدول ارتباط بین جفت‌های نمونه‌ها را با توجه به برچسب‌های واقعی Ground Truth و خوشه‌بندی انجام‌شده بررسی می‌کند. چهار حالت ممکن برای هر جفت نمونه در جدول:

**True Positive (TP)**: جفت نمونه‌ها در داده‌های واقعی و خوشه‌بندی یکسان هستند.

**False Positive (FP)**: جفت نمونه‌ها در خوشه‌بندی یکسان هستند اما در داده‌های واقعی متفاوت‌اند.

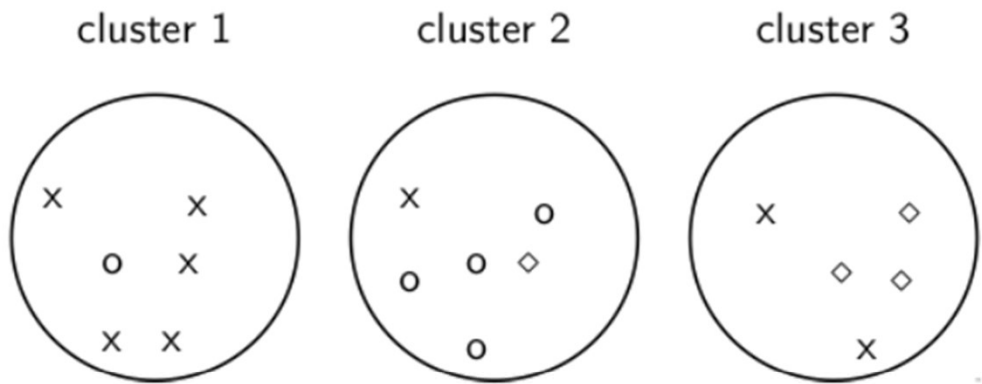
**False Negative (FN)**: جفت نمونه‌ها در داده‌های واقعی یکسان هستند اما در خوشه‌بندی متفاوت‌اند.

**True Negative (TN)**: جفت نمونه‌ها در داده‌های واقعی و خوشه‌بندی متفاوت‌اند.

	same cluster	different clusters
same class	true positives (TP)	false negatives (FN)
different classes	false positives (FP)	true negatives (TN)

- $TP+FN+FP+TN$  is the total number of pairs:  $\binom{n}{2}$  for  $n$  samples.

Definition:  $RI = \frac{TP+TN}{TP+TN+FP+FN}$



	same cluster	different clusters
same class	TP = 20	FN = 24
different classes	FP = 20	TN = 72

$RI = (20 + 72)/(20 + 20 + 24 + 72) \approx 0.68$



دانشگاه سمنان

دانشگاه سمنان

Semnan University

پردیس فرزانهگان

## ارزیابی داخلی: مجموع مربعات خطا Sum of Squared Error

SSE مجموع مربعات فاصله‌ها بین نقاط داده و مرکز خوشه Centroid مربوطه است. فرمول کلی SSE به صورت زیر است:

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x)$$

که در آن:

k: تعداد خوشه‌ها

C<sub>i</sub>: مجموعه نقاط در خوشه i

m<sub>i</sub>: مرکز خوشه i

x یک نمونه داده در خوشه C<sub>i</sub>

**SSE کمتر** نشان‌دهنده خوشه‌بندی بهتر است، زیرا نقاط داده به مرکز خوشه نزدیک‌تر هستند (تراکم بیشتر).

با افزایش تعداد خوشه‌ها، SSE معمولاً کاهش می‌یابد، زیرا خوشه‌های کوچک‌تر نقاط را دقیق‌تر توصیف می‌کنند.