

ملاحظة:

من الممكن أن يظهر جزء كود بايثون من اليمين إلى اليسار، لذا بفضل نسخة ولصقة في برنامجك المفضل لكتابة الكود أو عرضة بشكل أفضل في صفحات الملخصات في موقعي الشخصي:

<https://alioh.github.io/DSND-Notes-11/>

<https://alioh.github.io/DSND-Notes-12/>

<https://alioh.github.io/DSND-Notes-13/>

<https://alioh.github.io/DSND-Notes-14/>

<https://alioh.github.io/DSND-Notes-15/>

<https://www.alioh.com>

## Clustering - الدرس الأول - التجميع

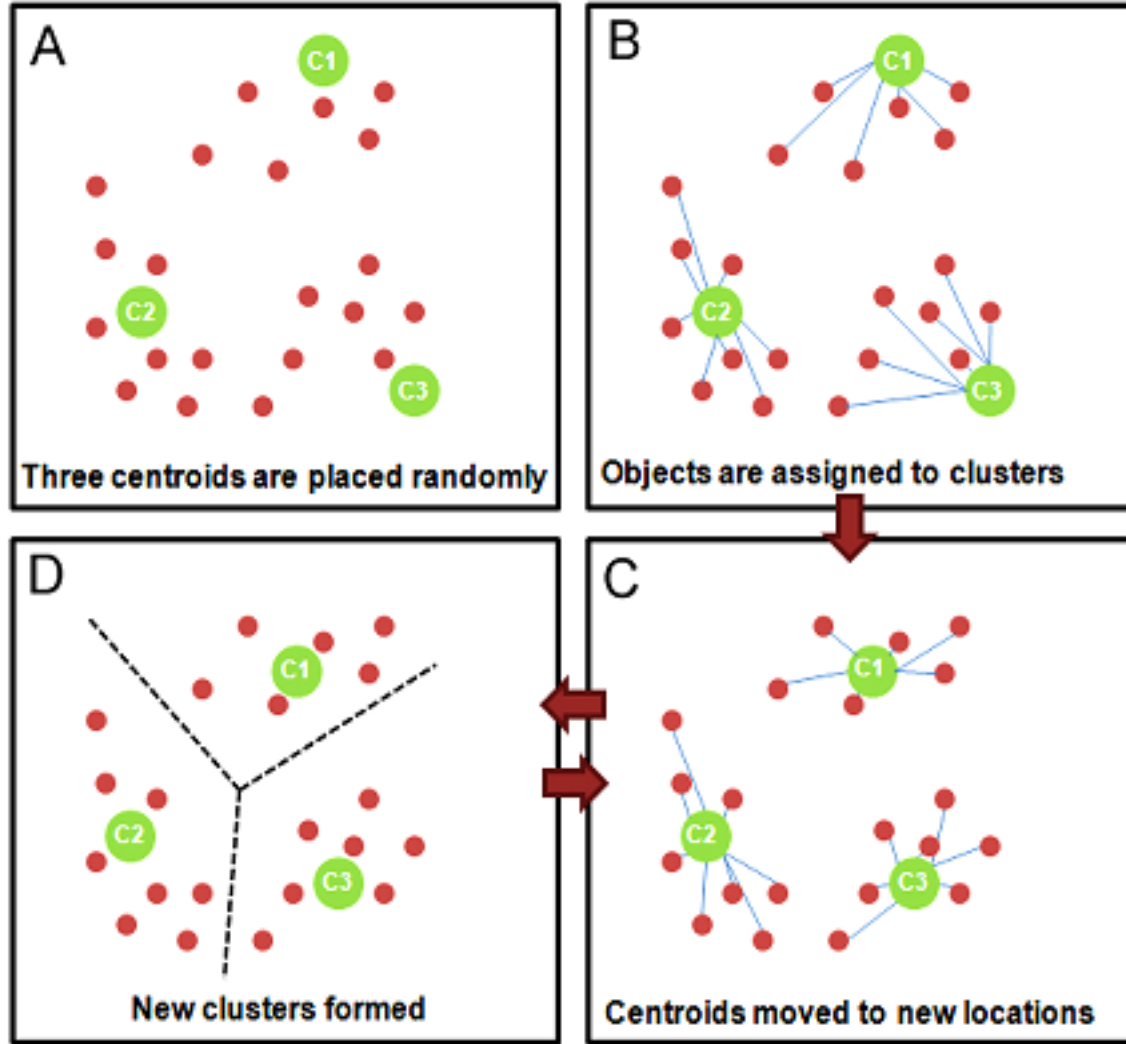
.عن التعلم غير الموجه بشكل عام، تعريفه وأمثلة عليه واشهر خوارزمياته [منشور سابق](#) كتبت في

### أنواع التعلم غير الموجه

- **Clustering** التجميع: جمع البيانات على شكل فئات حسب خصائص معينة تتشابه فيها.
- **Dimensionality reduction** تقليل الأبعاد: عرض البيانات بأقل عدد من الخصائص: **تقليص الأبعاد**.
- **Association Rules** ربط البيانات فيما بينها عن طريق محاولة اكتشاف علاقات بينها: **قواعد / قوانين الربط**.

### خوارزميات التعلم غير الموجه

#### تجميع بالمتوسطات K-Means



Netflix و Spotify تستخدم هذه الخوارزمية كثيراً في برامج تقديم الإقتراحات مثلاً في [هنا](#) تم شرحها سابقاً بشكل مفصل

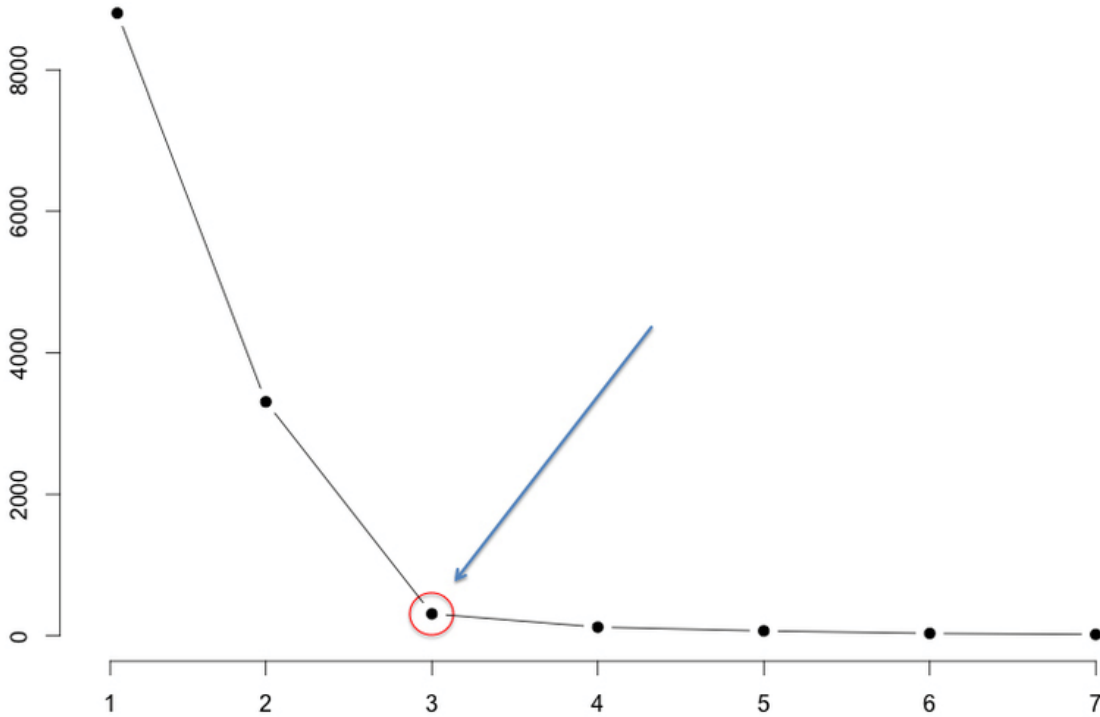
### K-Means طريقة عمل الخوارزمية

- Clusters. ويعني عدد مرات تقسيم (تجميع) البيانات إلى K يتم إعطاء الخوارزمية رقم
- Centroids. يتم رسم نقاط تعرف بالنقاط المركزية K بنفس الرقم

- تربط كل بيانات بما يقارب لها من نقاط مركزية.
- Cluster تقوم الخوارزمية بحساب المتوسط لجميع النقاط داخل كل
- تكرار الخطوتين الأخيرتين مرة أخرى حتى تصل لنقاط مركزية وتجميعة مناسبة لجميع النقاط

## Elbow Method

في البيانات لدينا. تستخدم هذه الطريقة Clusters تعني هنا عدد الـ K ، و K-Means في خوارزمية K طريقة لإيجاد عدد الـ :عندما لا نعرف عدد التصنيفات لدينا. طريقة عملها كالتالي



- نتبع الخطوات السابقة لعمل الخوارزمية
- Score يحتوي المودل على دالة Scikit-learn عند انشاء الخوارزمية بواسطة
- نظهر النتائج في رسم بياني يكون شكله كالتالي
- Cluster. في الصورة 3 هو العدد المناسب الذي يتوقف فيه التغير العالي في متوسط المسافة بين النقاط و الـ

$X = \text{Number of K} - Y = \text{Average Distance to Cluster}$

## مثال بايثون

Elbow Method وإستخدام الـ K-means طريقة كتابة خوارزمية

```
from sklearn.cluster import KMeans
```

```

import numpy as np
import matplotlib.pyplot as plt

# لنفرض لدينا بيانات معرفة في data
data = data

# نعرف المودل ونضيف له عدد ال Clusters
kmeans_4 = KMeans(n_clusters=4)

# نقوم بربط البيانات مع المودل
model_4 = kmeans_4.fit(data)

# للتوقع نقوم بالتالي
labels_4 = model_4.predict(data)

# matplotlib كرسم بياني بإستخدام مكتبة Elbow method الجزء التالي لعرض
# ونظهر نتائجها Clusters هنا نجرب من 1 إلى 10
No_of_clusters = range(1, 10)

# Clusters ننشأ قائمة تحتوي على 10 مودلز من 1 إلى 10
kmeans = [KMeans(n_clusters=i) for i in No_of_clusters]

# ننشأ قائمة ثانية تحتوي على نتائج كل مودل من التي أنشئناها بالخطوة السابقة
score = [np.abs(kmeans[i].fit(data).score(data)) for i in range(len(kmeans))]

# عرض النتائج في رسم بياني
plt.plot(No_of_clusters, score, linestyle='--', marker='o', color='b');
plt.xlabel('Number of Clusters')
plt.ylabel('Score')
plt.title('Elbow Curve')
plt.show()

```

## Feature Scaling

طريقة لإعادة تعيين قيمة النتائج حتى نحصل على بيانات ونتائج صحيحة، توجد طريقتين لذلك

Values	Normalized	Standardized
47	0.9302	1.1560
7	0.0000	-1.9267
21	0.3256	-0.8478
28	0.4884	-0.3083
41	0.7907	0.6936
49	0.9767	1.3102
50	1.0000	1.3872
25	0.4186	-0.5395
25	0.4186	-0.5395
35	0.6512	0.2312
24	0.3953	-0.6165

- Standardizing: Variance كل النتائج يساوي 0 وتباين Mean تحويل القيم إلى رقم جديد يكون متوسط 1. تكون النتائج عادة متقاربة من -3 إلى 3 مثلاً
- Normalizing: تحويل القيم إلى رقم جديد من 0 إلى 1

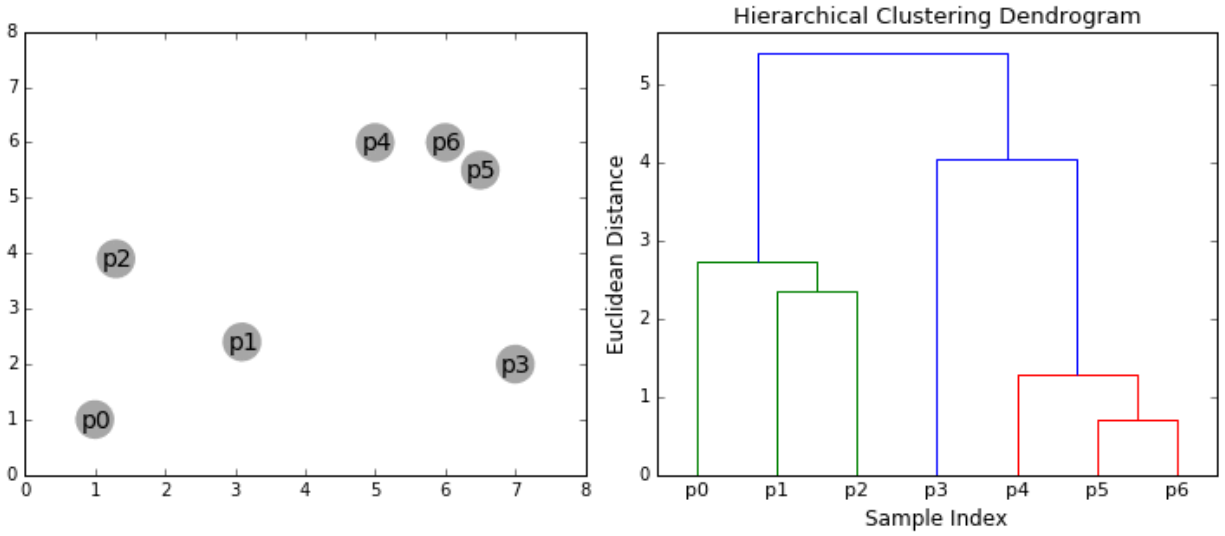
مثال لنتائج قبل وبعد التحويل لإحدى الطريقتين

-----

## - الدرس الثاني - التجميع القائم على الكثافة والهرمية Hierarchical and Density Based Clustering

### Hierarchical Clustering التجميع الهرمي

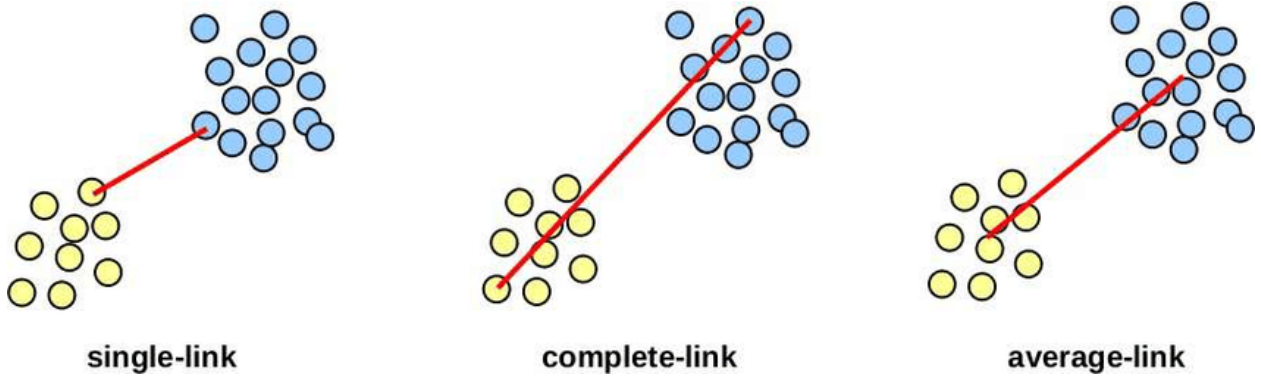
عن التجميع الهرمي. وشرحت طريقتها كالتالي [سابقاً](#) كتبت



- لكل نقطة لدينا Cluster التجميع الهرمي يُكوّن مجموعة أو.
- يتم جمع أقرب نقطتين مع بعضهم البعض وتُكوّن مجموعة جديدة.
- Cluster تبحث الخوارزمية عن نقطة أخرى قريبة لها وتضمها إلى نفس المجموعة.
- تُجمع فيها كل النقاط Cluster تستمر الخوارزمية بعمل الخطوتين السابقتين إلى أن تبقى مجموعة واحدة.

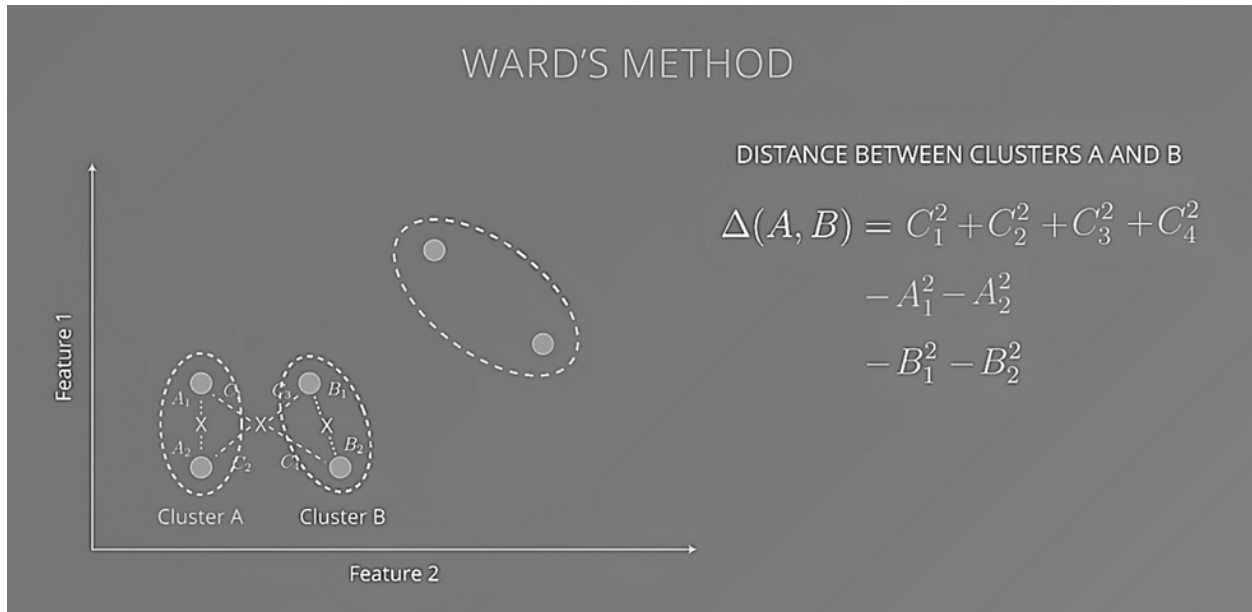
### أنواع التجميع الهرمي

في التجميع الهرمي Clusters توجد أربع أنواع لتحديد المسافات بين الـ



- **Single Link:** نقطتين لبعضها البعض في كل أقرب وأخرى عن طريق حساب Cluster تحدد المسافة بين كل Cluster.
- **Complete Link:** نقطتين لبعضها البعض في أبعد وأخرى عن طريق حساب Cluster تحدد المسافة بين كل Cluster.
- **Average Link:** في متوسط المسافة بين كل النقاط وأخرى عن طريق حساب Cluster تحدد المسافة بين كل Cluster.
- **Ward's Method:** بين Centroid تقوم بتكوين نقطة مركزية Ward's على عكس الطرق السابقة، طريقة ، بعد ذلك تحسب المسافة من النقاط إلى النقطة الجديدة التي أوجدتها سابقاً، تربيع المسافة Clusters مجموعتين

السابقة لكل نقطة ثم تُجمع جميعها. بعد ذلك نقوم بتربيع ثم طرح المسافة بين كل نقطة والنقطة المركزية في الCluster الخاص بها.



## كود بايثون

تطبيق عملي في بايثون للتجميع الهرمي

```
from sklearn.cluster import AgglomerativeClustering
from sklearn.metrics import adjusted_rand_score
from sklearn import preprocessing
from scipy.cluster.hierarchy import linkage
from scipy.cluster.hierarchy import dendrogram
import matplotlib.pyplot as plt
from sklearn import datasets
```

# Scikit-learn تحميل البيانات، في هذه المثال استخدمنا بيانات مرفقة من

# Iris بأسم

```
iris = datasets.load_iris()
```

# Clusters ننشأ المودل ونعطية عدد الـ

# Ward's Method النوع التلقائي المستخدم من المكتبة هي

# إذا اردنا إختيار نوع مختلف من الأربعة التي ذكرت سابقاً

# linkage نضيف الخيار

# ونعطية أحد الخيارات التالية:

# "ward", "complete", "average", "single"

# مثال:

```
# AgglomerativeClustering(n_clusters=3, linkage='complete')
```

```
ward = AgglomerativeClustering(n_clusters=3)
```

# fit\_predict الربط والتوقع بدالة واحدة

```
ward_pred = ward.fit_predict(iris.data)
```

```

# لتقييم المودل ومعرفة دقته
ward_ar_score = adjusted_rand_score(iris.target, ward_pred)
# النتيجة: 0.73119

# لتحسين النتائج أو إذا كانت لدينا فروقات كثيرة بين البيانات من
# Normalizing الممكن أن بال
# Normalization راجع الجزء 11 لشرح عن
normalized_X = preprocessing.normalize(iris.data)
# النتيجة: 0.88569

# لعرضها بشكل رسم بياني نقوم بالتالي
linkage_type = 'ward'
linkage_matrix = linkage(normalized_X, linkage_type)
plt.figure(figsize=(22,18))
dendrogram(linkage_matrix)
plt.show()

```

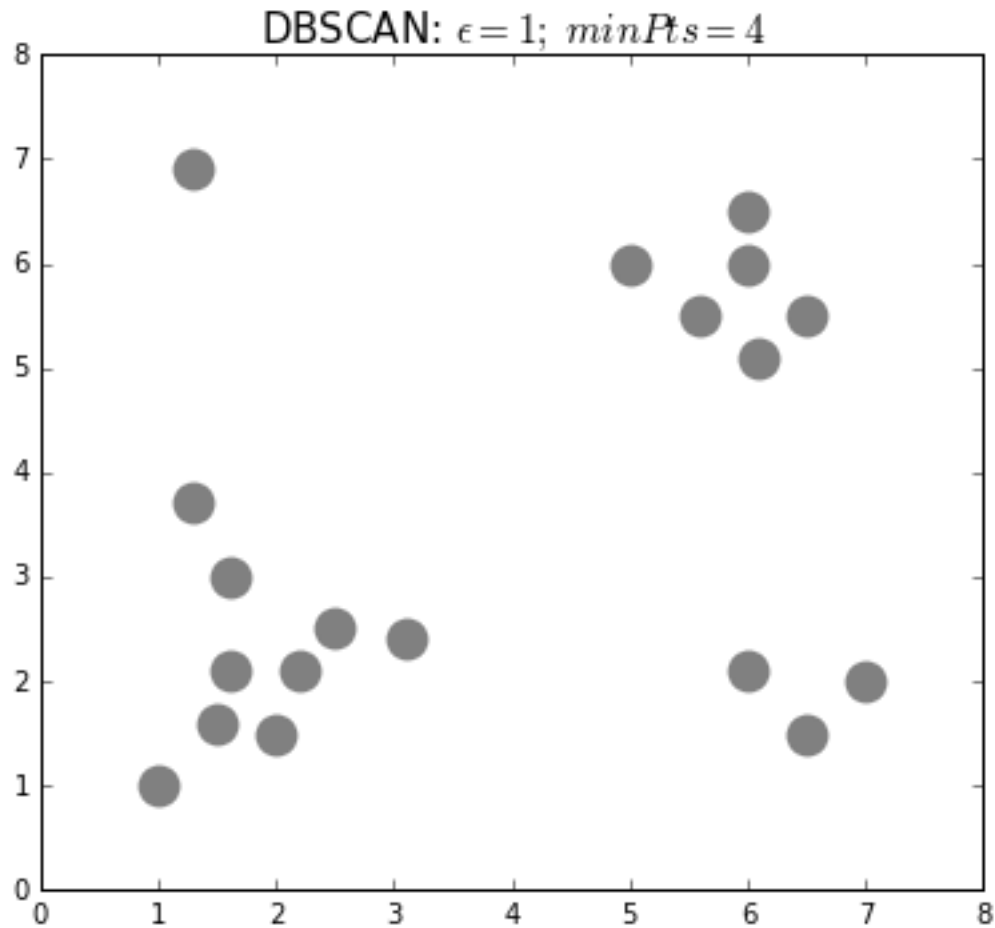
## Density-Based Clustering التجميع القائم على الكثافة

قبل انشاء المودل نعطي بعض المتغيرات

1. min\_samples أو Cluster عدد النقاط في كل
2. Epsilon  $\epsilon$  المسافة بين كل نقطة وأخرى

ع: نبدأ من نقطة معينة ونرى إذا كان هناك نقاط أخرى في المسافة التي حددناها





1. إذا لم يكن هناك نقاط أخرى، تحسب هذه النقطة كنقطة مزعجة ويتم تجاهلها.
2. إذا كان هناك نقاط أخرى ومجموعها أقل من الحد الأدنى لعدد النقاط الذي حددناه، نتجاهل هذه النقاط وننتقل لأي نقطة قريبة أخرى.
3. هنا، في حال وجدت Cluster إذا وجد نقاط أخرى ومجموع النقاط أعلى أو يساوي العدد الذي حددناه، نقوم بإنشاء لهذه النقطة، نجعلها جميعاً في Cluster نقاط أخرى توافي نفس الشروط بالحد الأدنى والمسافة وفي نفس الـ Cluster واحد.

## مثال بايثون

مثال عملي لطريقة إنشاء مودل للتجميع القائم على الكثافة

```
from sklearn.cluster import DBSCAN
```

```
# لتكن لدينا بيانات معرفة في المتغير data
data = data
```

```
# ننشأ المودل ونحدد المسافة eps
# والحد الأدنى من النقاط min_samples
```

```
dbscan = DBSCAN(eps=3, min_samples=2)
```

```
# fit_predict الربط والتوقع بدالة واحدة
```

```
# النتيجة 1- تعني ان القيمة مزعجه وتم تجاهلها
```

```
clustering_labels_1 = dbscan.fit_predict(data)
```

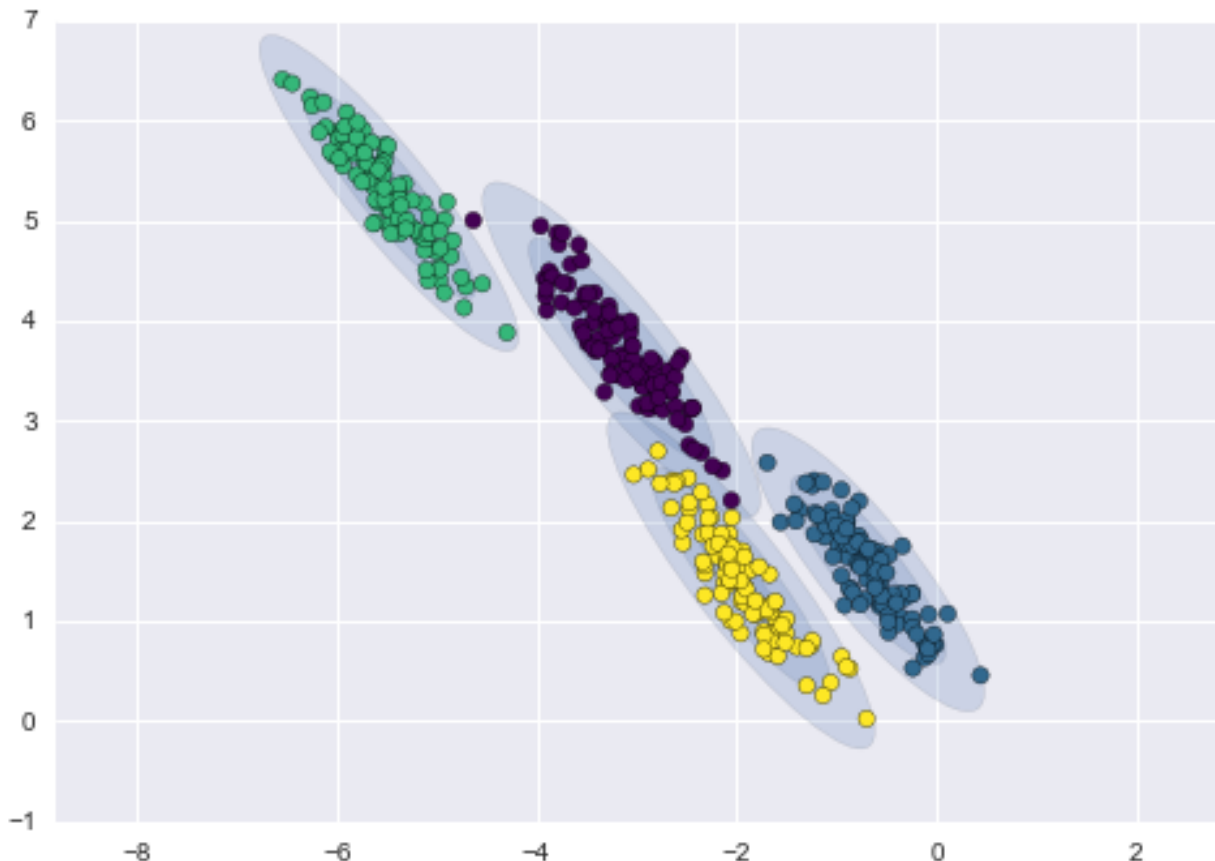
-----

## Gaussian Mixture Models and Cluster Validation - الدرس الثالث

### Gaussian Mixture Models (GMM)

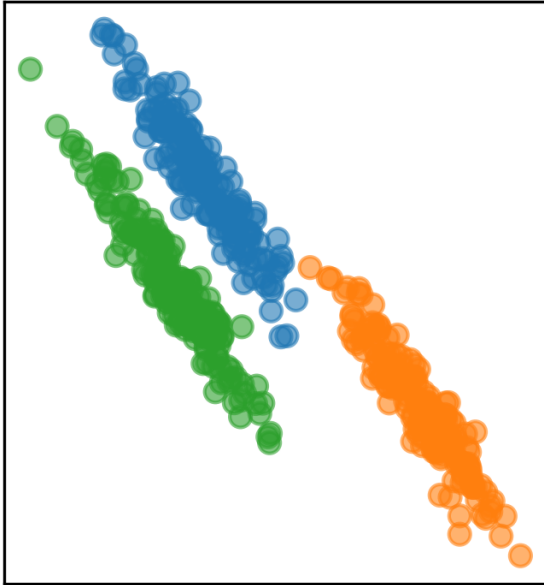
وهي طريقة لرسم Gaussian ولكن تأخذ بعين الاعتبار التباين. وسنستخدم فيها الـ K-mean خوارزمية تجميع، مشابهة لـ البيانات وفهم الخوارزمية لكيفية تشكيل البيانات لدينا.

- **Mean المتوسط** Cluster. وهو دائماً يكون في وسط الـ
- **Variance التباين** Cluster. ويشكل طريقة توسع البيانات في الـ Clusters ، وكما يظهر لنا طريقة توسعها بسبب عمل الخوارزمية على المتوسط والتباين لتشكيل الـ Clusters ، التي تؤدي عملها افضل على البيانات ذات الشكل الدائري K-mean على عكس ما كانت عليه خوارزمية

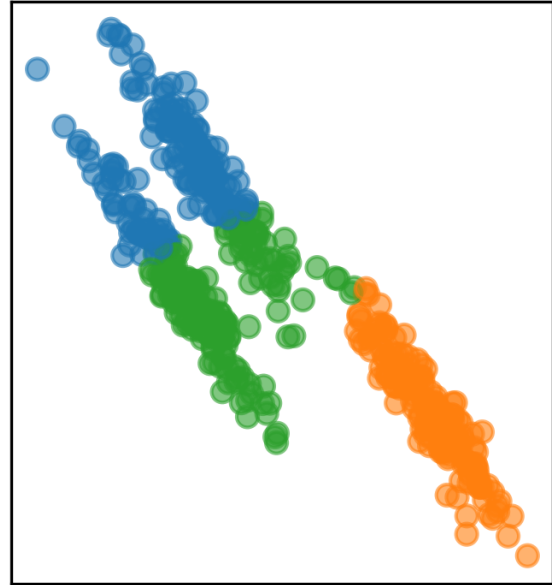


GMM و K-mean في Clusters مثال آخر للفرق في تكوين الـ

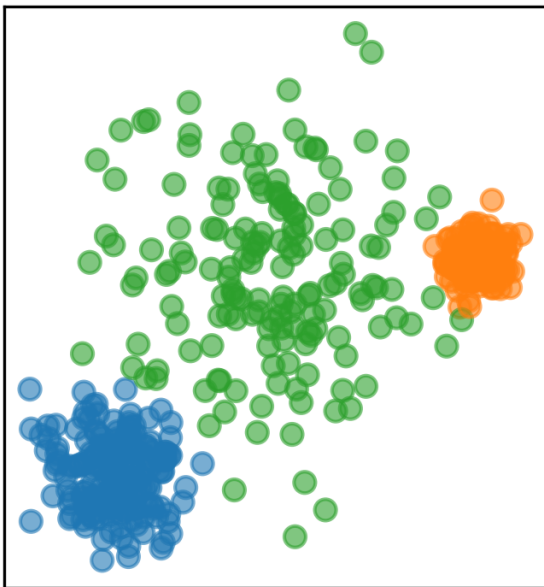
GaussianMixture



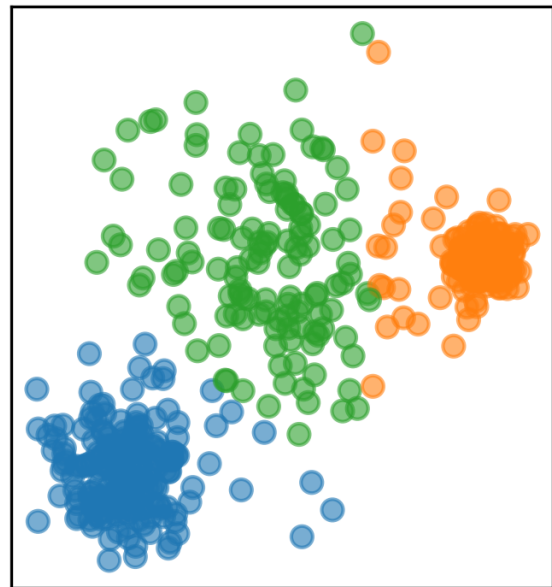
KMeans



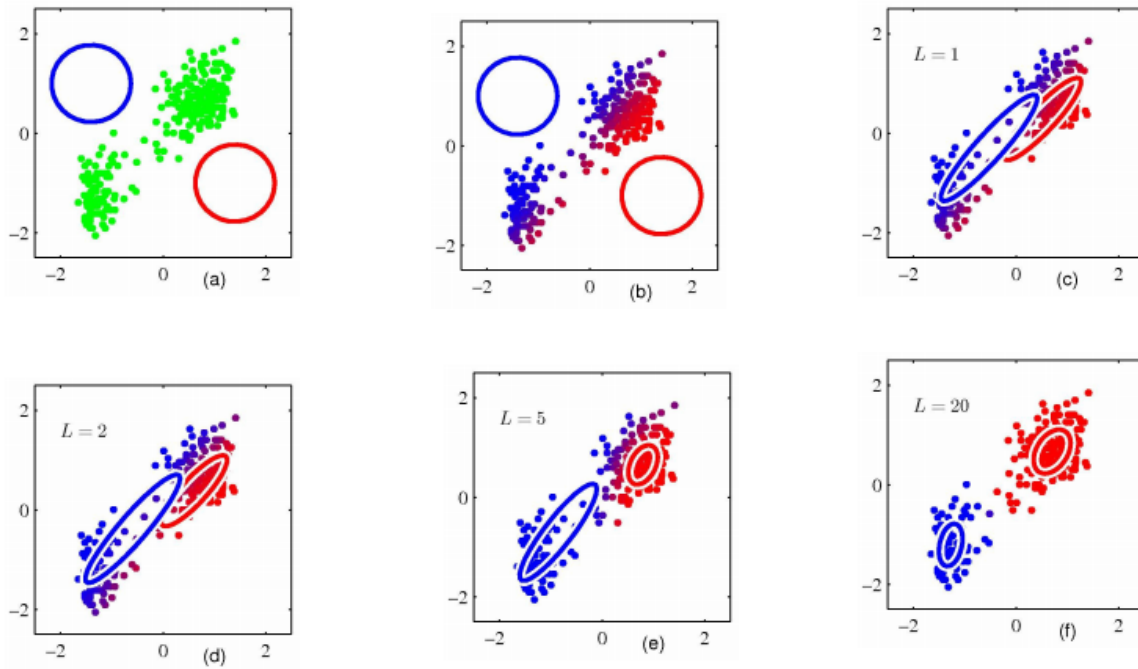
GaussianMixture



KMeans



**Expectation Maximization (EM)**



وطريقتها كالتالي GMM طريقة تتبع في خوارزمية

- وأختيار نقاط بداية عشوائية Clusters كعدد لـ  $K$  اختيار
- خاص بها Cluster نحدد لكل نقطة في البيانات
- نحسب المسافة والتباين بين كل نقطة والنقطة المركزية، ثم نقوم بتغيير موقع النقطة المركزية، ثم إعادة الحساب للمتوسط والتباين
- إذا توقف التغير في المتوسط يتم إختيار المنطقة كنقطة مركزية جديدة

## Clusters طريقة تحليل الـ

- **Feature Selecting:** اختيار المتغيرات او الخصائص المناسبة وتوجد طرق لإختيار الأفضل
- **Clustering Algorithm:** اختيار الخوارزمية المناسبة في التجميع وإعادة ضبطها حتى نصل لنتائج مرضية
- **Clustering Validation:** التحقق وعرض النتائج وتقييمها

## Cluster Validation

، في التعلم غير الموجه توجد ثلاث Precision و Accuracy في التعلم الموجه كانت لدينا خيارات تقييم للمودل كـ  
وهي Cluster خيارات لتقييم الـ

- **External Index:** المعروفة مسبقاً (Labels) مع النتائج الصحيحة Cluster يتم فيه مقارنة النتائج في الـ
- **Relative Index:** Clusters تقارن بين اثنين من الـ
- **Internal Index:** بدون الاعتماد على النتائج الصحيحة (ان لم تكن موجودة Cluster يتحقق من النتائج للـ (مثلاً).

$$S_i = \frac{b_i - a_i}{\text{Max}(b_i - a_i)}$$

**Silhouette Coefficient:** والمعادلة الخاصة بها (Labels) اذا لم تكن هناك نتائج Cluster طريقة لحساب اداء الـ كالتالي:

1.  $a$  = Cluster معدل المسافة من نقطة إلى النقاط الأخرى في نفس الـ.
2.  $b$  = Cluster معدل المسافة من نقطة إلى النقاط الأخرى في أقرب الـ. العملية على كل نقطة لدينا، ثم نجمع النتائج ونأخذ المعدل. وهو بالغالب رقم من 1 إلى 1

## مثال بايثون

GMM مثال لطريقة إنشاء مودل بخوارزمية

```
from sklearn.mixture import GaussianMixture
```

```
data = data
```

```
# Cluster 3 ننشأ المودل ونحدد عدد الـ
```

```
gmm = GaussianMixture(n_components=3)
```

```
# الربط والتوقع
```

```
gmm = gmm.fit(X)
```

```
pred_gmm = gmm.predict(X)
```

-----

## الدرس الرابع - Dimensionality Reduction and PCA

تقوم الخوارزمية بجمع البيانات ومحاولة إيجاد انماط فيما Patterns. في التعلم غير الموجه لإكتشاف الأنماط PCA تستخدم خاصة بها بجانب ما يشابهها. الفيديو التالي يشرح طريقة عملها Cluster بينها، ثم جمع كل واحد في

### Latent Features - Feature Selection

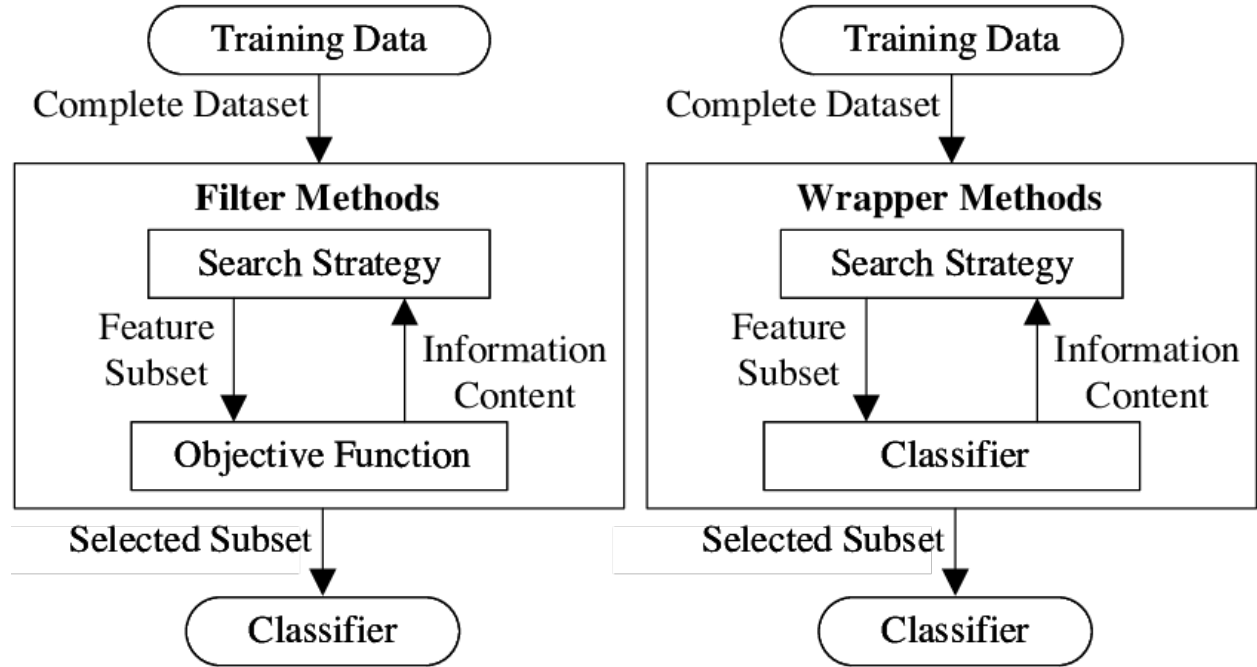
:متغيرات ليست موجودة في البيانات الأصلية لدينا، لنقل أن لدينا بيانات منازل وفيها المتغيرات التالية:

- حجم الأرض
- عدد الغرف

- مساحة المنزل
- مساحة المخزن
- عدد غرف النوم
- عدد دورات المياة
- معدل الجرائم في الحي
- عدد المدارس في الحي
- الضرائب في الحي
- متوسط مصادر الدخل للجيران
- متوسط جودة الهواء
- المسافة إلى الخط السريع

في المتغيرات السابقة ممكن ان نستخرج متغيرين جدد غير الموجودة لدينا وهي لأول 6 متغيرات وآخر 6 متغيرات، بذلك قللنا عدد المتغيرات لدينا في البيانات معلومات الحي والثاني حجم المنزل المتغيرات الجدد هي

توجد عدة طرق لإختيار المتغيرات وهي:



- **Filter methods:** ومن مسماها ممكن ان نتوقع أنها تقوم بفلتره المتغيرات الموجودة وترتيبها حسب اهميتها [هنا](#) لدينا، يتم بعد ذلك تقييمها بعدة طرق (تفاصيل أكثر عن هذه الطرق
- **Wrapper Methods:** هذه الطريقة تقوم بالتحقق من اهمية المتغير مباشرة مع المودل، اذا انتجت المودل افضل نتيجة له فسيتم إختيار تلك المتغيرات

## مثال بايثون

PCA مثال عملي طويل ومفصل لطريقة عمل خوارزمية [هنا](#): في الرابط التالي

## PCA متى نستخدم

لدينا. في بعض الأحيان في الصور مثلاً، قللنا Features تستخدم الخوارزمية متى ما أردنا التقليل من عدد المتغيرات المتغيرات من 700 متغير إلى 30 متغير وحصلنا على نفس النتائج.

## Random Projection & ICA - الدرس الخامس

أو المتغيرات لدينا. لنفترض أن Dimensions لتقليل المتغيرات. هدف الطريقة الرئيسي هو تقليل الـ PCA هي حل بديل لـ Random Projection وتشغيل دالة Scikit-learn لدينا قاعدة بيانات فيها 12000 عامود أو متغير، عند إدخالها في. عليها ستعود لنا البيانات بعدد عوامل حوالي 6000.

### مثال بايثون

عمل Random Projection في بايثون و Scikit-learn

```
from sklearn import random_projection

data = data

# إنشاء المودل الخاص ب Random Projection
# عدد المتغيرات المطلوب n_components من الممكن تحديد عدد
# eps وتحديد قيمة
# eps = 0.1 و n_components = 'auto' اذا تركت بدون تحديد فتكون قيمة
# وهي القيم التلقائية
rp = random_projection.SparseRandomProjection()

# تشغيل المودل على البيانات لدينا
new_data = rp.fit_transform(data)
```

### Independent Component Analysis - ICA

مثال عليها في ملفات الصوت، اذا كان لدينا ثلاث ادوات صوتية Random Projection و PCA طريقة أخرى مشابهة لـ. تعمل بنفس الوقت، تشغيل هذه الخوارزمية عليها ستحاول فصل كل اداة لوحدها.

### مثال بايثون

عمل Random Projection في بايثون و Scikit-learn

```
from sklearn.decomposition import FastICA

# البيانات لدينا متكونة من قائمة تحتوي على 3 ملفات صوتية
data = list(zip(audio_1, audio_2, audio_3))

# إنشاء المودل الخاص
# عدد المتغيرات المطلوب n_components تحديد عدد
```

```
ica = FastICA(n_components=3)
```

```
# تشغيل المودل على البيانات لدينا  
new_data = ica.fit_transform(data)
```

أنتهى المخلص

-----

الملخصات منشورة بشكل مفصل وبترتيب أفضل في صفحتي الشخصية هنا

<https://alioh.github.io/DSND-Notes-11/>

<https://alioh.github.io/DSND-Notes-12/>

<https://alioh.github.io/DSND-Notes-13/>

<https://alioh.github.io/DSND-Notes-14/>

<https://alioh.github.io/DSND-Notes-15/>

<https://www.alioh.com>

-----

مصادر:

1. [mubaris](#) - K-Means Clustering
2. [Quora](#) - Feature Scaling
3. [saedsayad](#) - Hierarchical Clustering
4. [dashee87](#) - Hierarchical Clustering 2
5. [jakevdp](#) - Gaussian Mixture
6. [amueller](#) - Clustering and Mixture Models
7. [columbia.edu](#) - Gaussian Mixture Models
8. [kent.edu](#) - Cluster Validation
9. [towardsdatascience](#) - Clustering Analysis
10. [analyticsvidhya](#) - Feature Selection methods