# (Clustering Algorithms) مقارنة خوارزميات التجميع

# جدول المقارنة الرئيسي

الخاصية	K- Means	DBSCAN	Hierarchical	Gaussian Mixture	Spectral Clustering	MeanShift	OPTICS
التعقيد الحسابي	O(nkt)	O(n log n)	O(n³)	O(nkd²)	O(n³)	O(n²)	O(n²)
حاجة لتحديد عدد المجموعات	نعم	И	لا (مرن)	نعم	نعم	И	Л
الأداء مع البيانات الكبيرة	ممتاز	جيد	ضعيف جداً	متوسط	ضعیف	ضعیف	متوسط
شكل المجموعات	دائري	أي شكل	أي شكل	إهليلجي	معقد	أي شكل	أي شكل
مقاومة الضوضاء	ضعيفة	ممتازة	متوسطة	جيدة	متوسطة	جيدة	ممتازة
مقاومة الـ Outliers	ضعيفة جداً	ممتازة	ضعيفة	جيدة	متوسطة	جيدة	ممتازة
حساسية للمعايرة	متوسطة	عالية	منخفضة	متوسطة	عالية	عالية	متوسطة
حساسية لاختيار البداية	عالية	И	И	متوسطة	И	И	И
إنتاج مجموعات متوازنة الحجم	نعم	И	И	نعم	متوسط	И	И
التعامل مع الأبعاد العالية	متوسط	ضعیف	ضعیف	ضعیف	جيد	ضعیف	ضعیف
سهولة التفسير	عالية جداً	عالية	ممتازة	متوسطة	متوسطة	عالية	عالية

# التفاصيل الفنية

### 1. **©** K-Means

- مجموعات بحيث كل نقطة تنتمي للمركز الأقرب K المبدأ: يقسم البيانات لـ •
- :الخوارزمية
  - مراكز عشوائية K اختيار .1
  - تخصيص كل نقطة للمركز الأقرب .2
  - إعادة حساب المراكز .3
  - تكرار حتى الاستقرار .4

### :المعاملات المهمة

- (n\_clusters): عدد المجموعات (يجب تحديده)
- (init): طريقة اختيار المراكز (k-means++')
- عدد المحاولات (10 افتراضي) : n\_init
- (300) تكرارات (300) •

### 2. Q DBSCAN (Density-Based Spatial Clustering)

- المبدأ: يجمع النقاط القريبة عالية الكثافة، يتجاهل الضوضاء •
- المفاهيم: Core points, Border points, Noise points
- :المعاملات المهمة
  - المسافة القصوى بين النقاط (حساس جداً!) (eps
  - (min\_samples): (10-5) أقل عدد نقاط لتكوين مجموعة
- القوة: يكتشف أي شكل، يتعامل مع الضوضاء •
- بطيء مع البيانات عالية الأبعاد ،eps الضعف: صعب اختيار •

### 3. Hierarchical Clustering

- نوعان (من أعلى لأسفل)، Divisive (من أسفل لأعلى)، كانتناب
- Agglomerative: کل نقطة مجموعة o یدمج الأقرب o حتى مجموعة واحدة
- المخرجات: Dendrogram المخرجات
- :المعاملات المهمة
  - عدد المجموعات النهائي :(n\_clusters)
  - (linkage): طريقة حساب المسافة ('ward', 'complete', 'average')
- مسبقاً، سهل التصور K القوة: لا يحتاج تحديد •
- O(n³) الضعف: بطيء جداً •

### 4. P Gaussian Mixture Models (GMM)

المبدأ: يفترض أن البيانات مزيج من توزيعات غاوسية •

- (soft clustering) يعطى: احتمالية انتماء كل نقطة لكل مجموعة
- :المعاملات المهمة
  - (n\_components): عدد الـ Gaussians
  - (covariance\_type): شكل التوزيع ('full', 'tied', 'diag', 'spherical')
- uncertainties القوة: مرونة في الأشكال، يعطى
- الضعف: يفترض التوزيع الغاوسي •

### 

- للتجميع eigenvalues ويستخدم graph المبدأ: يحول البيانات لـ
- القوة: يتعامل مع الأشكال المعقدة والمنحنيات
- :المعاملات المهمة
  - (n\_clusters): عدد المجموعات
  - (affinity): طريقة بناء الـ graph ('rbf', 'nearest\_neighbors')
- صعب اختيار المعاملات ،O(n³)، الضعف: بطيء •

### 6. MeanShift

- المبدأ: يجد مراكز الكثافة العالية عبر تحريك النقاط نحو المتوسط المحلي •
- القوة: يحدد عدد المجموعات تلقائياً، يتعامل مع أي شكل
- :المعاملات المهمة
  - عرض النافذة (حساس جداً!) :(bandwidth
- الضعف: بطيء جداً، حساس لاختيار

### 7. 🔬 OPTICS (Ordering Points To Identify Clustering Structure)

- ينتج ترتيب يوضح بنية التجميع ،DBSCAN المبدأ: تطوير لـ •
- DBSCAN القوة: أقل حساسية للمعاملات من •
- :المعاملات المهمة
  - (min\_samples): أقل عدد نقاط
  - أقصى مسافة للبحث :(max\_eps
- المخرجات: Reachability plot

### 8. 😕 Birch (Balanced Iterative Reducing and Clustering)

- للتجميع الهرمي (CF (Clustering Features **المبد**أ: يبني شجرة •
- القوة: سريع جداً، مناسب للبيانات الكبيرة جداً
- :المعاملات المهمة

- (n\_clusters): عدد المجموعات
- (threshold): عتبة الـ radius
- الضعف: يفترض شكل دائري للمجموعات

### 9. MiniBatch K-Means

- تستخدم عينات صغيرة K-Means المبدأ: نسخة محسنة من •
- العادي K-Means القوة: أسرع بكثير من •
- :المعاملات المهمة
  - (batch\_size): حجم الـ batch (100-1000)
- الاستخدام: مع البيانات الضخمة جداً •

### 10. Propagation

- exemplars المبدأ: كل نقطة ترسل رسائل للأخرى لتحديد الـ
- القوة: لا يحتاج تحديد عدد المجموعات مسبقاً
- :المعاملات المهمة
  - (preference): تفضيل النقاط لتكون exemplars
  - (damping): (1.0-0.5) معامل التخميد
- O(n²) الضعف: بطيء جداً

### متى نستخدم كل خوارزمية؟

:حسب شكل المجموعات

### :مجموعات دائرية/كروية

- K-Means (الأسرع والأبسط)
- MiniBatch K-Means (للبيانات الكبيرة)
- Birch (للبيانات الضخمة)

### :أشكال معقدة ومنحنيات

- DBSCAN (الأفضل عموماً)
- Spectral Clustering
- MeanShift

### :مجموعات إهليلجية

Gaussian Mixture Models

## وحسب حجم البيانات

حجم البيانات	الأفضل	البديل
< 1K	أي خوارزمية	للتصور Hierarchical
1K-10K	K-Means، DBSCAN	GMM
10K-100K	K-Means، DBSCAN	Spectral
100K-1M	MiniBatch K-Means	Birch
> 1M	Birch, MiniBatch K-Means	-

### دسب معرفة عدد المجموعات 🔎

### :أعرف عدد المجموعات

- **K-Means** (الأفضل)
- Spectral Clustering
- GMM

### :لا أعرف عدد المجموعات

- **DBSCAN** (يحدد تلقائياً)
- MeanShift
- Affinity Propagation
- Hierarchical (مرن)

### مسب وجود ضوضاء 🛕

### بيانات نظيفة:

- K-Means
- GMM

### :بيانات بها ضوضاء

- **DBSCAN** (الأمثل)
- OPTICS
- MeanShift

### استراتيجيات للاختيار العملي



### وللمبتدئين 🚀

(بسيط ومفهوم) **K-Means** ابدأ بـ .1

- 2. إذا النتائج مش منطقية → DBSCAN
- 3. للتصور والفهم → Hierarchical

### وللمحترفين 6

- تحليل البيانات أولاً (الشكل، الحجم، الضوضاء) .1
- 2. **K-Means** ≤ baseline
- للأشكال المعقدة **DBSCAN**
- للنتائج الاحتمالية 4. GMM

### ونصائح ذهبية

- لاختيار K-Means في K لاختيار K في K
- استخدم جودة التجميع Silhouette Analysis
- تصور البيانات قبل اختيار الخوارزمية •
- جرب **عدة خوارزميات** وقارن النتائج •

### M Recordsتطبيق عملي لـ 2

### :للبيانات بحجم 2 مليون

### الخيارات العملية 🔽

- 1. Birch: الأسرع والأمثل للحجم الضخم
- 2. MiniBatch K-Means: سريع ودقيق
- العادي: إذا الذاكرة كافية 3. **K-Means**

### : تجنب تماماً

- Hierarchical (اسیأخذ أیام!)
- Affinity Propagation (مش عملي)
- MeanShift (بطیء جداً

### هقترح Pipeline 🍪

تقییم النتائجo (PCA/UMAP) o Birch/MiniBatch K-Means o تقلیل الأبعاد

### التحسين الأداء

- لتقليل الأبعاد أولاً PCA استخدم •
- جرب sampling جرب

• استخدم parallel processing استخدم

### مقاييس التقييم

المقياس	الاستخدام	المعنى
Silhouette Score	جودة التجميع العامة	إلى 1 (أعلى أفضل) 1-
Calinski-Harabasz	وضوح الفصل	أعلى أفضل
Davies-Bouldin	تماسك المجموعات	أقل أفضل
Inertia	فقط K-Means للـ	(Elbow method) أقل أفضل

### نصائح التطبيق العملي

### قبل التجميع 🦴

- نظف البيانات من القيم المفقودة •
- طبع البيانات (StandardScaler)
- قلل الأبعاد إذا كانت عالية (>50) •

### وأثناء التجميع 🐑

- جرب قيم مختلفة للمعاملات •
- استخدم random\_state التخدم
- **صور النتائج** للفهم البصري •

### بعد التجميع 📊

- **فحص المجموعات** منطقياً •
- حلل خصائص كل مجموعة •
- **تحقق من التوزيع** (متوازن؟) •

الخلاصة: لمعظم المشاكل، ابدأ بـ **K-Means** وللبيانات الضخمة استخدم **K-Means**!