Clustering للـ Internal Validation Metrics

- نظرة عامة على المقاييس 🎯
- 1. Silhouette Score 🙎
- 2. Calinski-Harabasz Index 📊
- 3. Davies-Bouldin Index
- 4. Dunn Index 6

المقارنة الأساسية 📊

| المقياس | المدى | الأفضل | التعقيد الحسابي | الحساسية للضوضاء |
|-------------------|---------|------------|-----------------|------------------|
| Silhouette Score | [-1, 1] | أعلى (→ 1) | O(n²) | متوسطة |
| Calinski-Harabasz | [0, ∞) | أعلى | O(n) | منخفضة |
| Davies-Bouldin | [0, ∞) | أقل (→ 0) | O(n) | متوسطة |
| Dunn Index | [0, ∞) | أعلى | O(n²) | عالية |
| Dunn Index | [0, ∞) | اعلی | O(n²) | عالية |

التحليل التفصيلي 🔍

1. Silhouette Score

طريقة الحساب

i: لكل نقطة

a(i) = متوسط المسافة داخل الـ cluster

آخر cluster متوسط المسافة لأقرب = (b(i

s(i) = (b(i) - a(i)) / max(a(i), b(i))

s(i) متوسط کل = Silhouette Score

المزايا 🔽

- **سهل التفسير**: القيم واضحة ومفهومة •
- **يراعي التماسك والانفصال**: يجمع بين المعيارين •
- **مفصل**: يمكن حساب قيمة لكل نقطة ●
- **مستقر**: نتائج ثابتة نسبياً •

العيوب 🗶

- مع البيانات الكبيرة (O(n²) :**بطء الحساب**
- الكروية clusters حساس للشكل: يفضل الـ •
- متأثر بالكثافة: مشاكل مع الكثافات المختلفة •
- **حساس للأبعاد**: أداء ضعيف في الأبعاد العالية •

أفضل استخدام 🌀

- البيانات متوسطة الحجم البيانات متوسطة الحجم
- الكروية أو المتشابهة في الشكل clusters الـ •
- عندما تريد تحليل مفصل لكل نقطة •

2. Calinski-Harabasz Index (Variance Ratio Criterion)

طريقة الحساب

 $CH = [SS_B / (k-1)] / [SS_W / (n-k)]$

SS_B = Sum of squares بين الـ clusters

SS_W = Sum of squares داخل الـ clusters

k = عدد الـ clusters

عدد النقاط = n

المزايا 🔽

- سریع جداً: O(n) complexity
- outliers مقاوم للضوضاء: أقل تأثراً بالـ •
- Silhouette يعمل مع أشكال مختلفة: أكثر مرونة من
- **ذاكرة قليلة**: لا يحتاج

العيوب 🗶

- **صعب التفسير**: الأرقام المطلقة غير واضحة المعنى •
- over-clustering للـ endency أكبر clusters يفضل عدد
- **حساس للبيانات الخطية**: مشاكل مع البيانات غير الكروية •
- المتشابهة في الحجم clusters يفضل الـ :clusters متأثر بحجم الـ

أفضل استخدام 🎯

- البيانات الكبيرة (> 100K points)
- عندما السرعة مهمة •
- المتوازنة في الحجم clusters الـ
- كمقياس أولي سريع •

3. Davies-Bouldin Index

طريقة الحساب

```
لكل cluster i: S_i = Ji متوسط المسافة داخل الـ cluster M_i = M_i =
```

المزايا 🔽

- سریع: O(n) complexity
- بديهى: يقيس نسبة التماسك للانفصال
- لا يحابي عدد معين :**clusters مستقل عن عدد الـ**
- مقياس واضح: كلما قل كلما أفضل •

العيوب 🗶

- يفترض الشكل الكروي: مشاكل مع الأشكال المعقدة •
- clusters يعتمد على مراكز الـ :centroids حساس للـ
- متأثر بالحجم: مشاكل مع الأحجام المختلفة •
- المتداخلة clusters قد يكون مضلل: مع الـ

أفضل استخدام 🎯

- الكروية أو المتشابهة clusters الـ
- عندما تريد مقياس سريع وبسيط
- مختلف clusters للمقارنة بين عدد •
- البيانات متوسطة إلى كبيرة الحجم •

4. Dunn Index

طريقة الحساب

Dunn = min(d(i,j)) / max(d'(k))

d(i,j) =ا قل مسافة بين cluster i و d'(k) = أكبر مسافة داخل cluster k

المزايا 🔽

- مفهوم بديهى: يقيس الانفصال الواضح
- الضعيفة clusters يكتشف الـ outliers:
- مناسب للأشكال المختلفة: أكثر مرونة •
- تفسير واضح: كلما زاد كلما أفضل

العيوب 🗶

- أو أكثر (O(n²) :**بطء شديد**
- واحد يفسد كل شيء outlier :حساس جداً للضوضاء
- **غير مستقر**: نتائج متذبذبة •
- **صعب الحساب**: يحتاج كل المسافات •

أفضل استخدام 6

- البيانات الصغيرة (< 1K points)
- outliers البيانات النظيفة بدون
- عندما الانفصال الواضح مهم •
- كمقياس تكميلي وليس أساسي •

توصيات الاستخدام 🙎

(1K >) للبيانات الصغيرة

- الأكثر دقة وتفصيلاً 1. Silhouette Score
- كلتحقق من الانفصال 2. **Dunn Index**
- مقياس سريع إضافي 3. **Davies-Bouldin**

(1K - 10K) للبيانات المتوسطة

1. **Silhouette Score** - توازن جید

- 2. Calinski-Harabasz للسرعة
- مقياس ثانوي 3. **Davies-Bouldin**

(10K <) للبيانات الكبيرة

- الأسرع والأكثر عملية 1. Calinski-Harabasz
- مقياس ثانوي سريع 2. Davies-Bouldin
- على عينة من البيانات فقط 3. Silhouette

حسب نوع البيانات 📈

البيانات الكروية/المتجانسة

- الأول: Silhouette Score
- الثانى: Davies-Bouldin
- الثالث: Calinski-Harabasz

البيانات غير المنتظمة/المعقدة

- ا**لأول** (بحذر) Dunn Index
- الثاني: Silhouette Score
- الثالث: Calinski-Harabasz

Outliers البيانات مع

- الأول: Calinski-Harabasz
- الثانى: Davies-Bouldin
- تجنب: Dunn Index

النصائح العملية 🦴

Clusters عند اختيار عدد الـ

- للدقة Silhouette استخدم .1
- للسرعة Calinski-Harabasz استخدم .2
- **قارن النتائج** من عدة مقاييس .3

عند مقارنة الخوارزميات

سریع Calinski-Harabasz ابدأ بـ 1

- أكثر دقة **Silhouette تأكد بـ** 2.
- استخدم 3. Davies-Bouldin

للحصول على أفضل النتائج

- اجمع بين عدة مقاييس •
- اعتبر طبيعة بياناتك •
- الدرس توزيع القيم وليس فقط المتوسط
- اختبر على بيانات متنوعة •

الخلاصة النهائية 📊

| المقياس | أفضل لـ | تجنبه مع | السرعة | الدقة |
|-------------------|-------------------|----------------------|--------|-------|
| Silhouette | التحليل التفصيلي | البيانات الكبيرة | | **** |
| Calinski-Harabasz | البيانات الكبيرة | الأشكال المعقدة | * | *** |
| Davies-Bouldin | المقارنات السريعة | الأشكال غير الكروية | x° | *** |
| Dunn | الانفصال الواضح | Outliers البيانات مع | | ** |

:التوصية الذهبية

للحصول على Silhouette Score للاستكشاف السريع، ثم تأكد بـ Silhouette Score ابدأ بـ تقييم دقيق ومفصل.