مقارنة مقاييس المسافة في التجميع (Distance Metrics in Clustering)

جدول المقارنة الرئيسي

الخاصية	Euclidean	Manhattan	Chebyshev	Cosine	Jaccard	Hamming	Mahaland
التعقيد الحسابي	O(d)	O(d)	O(d)	O(d)	O(d)	O(d)	O(d²)
حساسية للأبعاد العالية	عالية جداً	متوسطة	منخفضة	منخفضة جدآ	متوسطة	منخفضة	عالية
حساسية للمعايرة	عالية جداً	عالية	متوسطة	منخفضة جداً	منخفضة	منخفضة	منخفضة
يتعامل مع الاتجاه	И	И	И	نعم	И	Л	Л
مناسب للبيانات المستمرة	ممتاز	جيد جداً	جيد	ممتاز	ضعیف	ضعیف	ممتاز
مناسب للبيانات الفئوية	ضعیف	ضعیف	ضعیف	ضعیف	ممتاز	ممتاز	ضعیف
مناسب للبيانات الثنائية	ضعیف	ضعیف	ضعیف	جيد	ممتاز	ممتاز	ضعیف
مقاومة الضوضاء	متوسطة	جيدة	جيدة جداً	جيدة	جيدة	جيدة جداً	ممتازة
سهولة التفسير	عالية جداً	عالية جداً	عالية	متوسطة	عالية	عالية جداً	متوسطة
مناسب للنصوص	ضعیف	ضعیف	ضعیف	ممتاز	جيد	ضعیف	ضعیف
يحافظ على الشكل الهندسي	نعم	نعم	نعم	И	-	-	نعم

التفاصيل الفنية

1. 🔪 Euclidean Distance (المسافة الإقليدية)

- المعادلة: $d = \sqrt{(\Sigma(x_i y_i)^2)}$
- "المعنى: المسافة المباشرة "كما يطير الطائر •
- التشبيه: المسافة بين نقطتين على الخريطة بخط مستقيم
- :الاستخدام

```
python
```

from scipy.spatial.distance import euclidean

dist = euclidean([1, 2], [4, 6]) # = 5.0

- أفضل مع: البيانات المعايرة، التوزيع الطبيعي •
- المشكلة: حساس جداً للأبعاد العالية (curse of dimensionality)

(مسافة المدينة) Manhattan Distance

- المعادلة: $(d = \Sigma | x_i y_i|)$
- المعنى: مجموع الفروقات المطلقة •
- التشبيه: المسافة في شوارع نيويورك (مربعات الشطرنج) •
- euclidean من outliers المميزات: أقل تأثراً بالـ
- الاستخدام

python

from scipy.spatial.distance import manhattan

dist = manhattan([1, 2], [4, 6]) # = 7

3. 1 Chebyshev Distance (مسافة الملك)

- المعادلة (d = max|x_i y_i|)
- المعنى: أكبر فرق في أي بُعد •
- التشبيه: حركة الملك في الشطرنج •
- الاستخدام: عندما البُعد الأسوأ هو المهم
- مثال: تقييم الأداء (أسوأ مقياس يحدد الجودة)

4. Cosine Distance/Similarity

- المعادلة: (cosine_sim = (A·B)/(||A|| × ||B||)
- Cosine Distance: 1 cosine_similarity

- المعنى: يقيس الزاوية بين المتجهين، يتجاهل الحجم
- القوة: مثالي للنصوص وعالية الأبعاد •
- مثال •

python

from sklearn.metrics.pairwise import cosine_similarity sim = cosine_similarity([[1, 2]], [[2, 4]]) # = 1.0 (نفس الاتجاه)

5. **@** Jaccard Distance/Similarity

- للبيانات الثنائية والمجموعات
- J = |A ∩ B| / |A ∪ B|
- **Jaccard Distance**: 1 Jaccard_similarity
- **الاستخدام**: مقارنة المجموعات، البيانات الثنائية •
- **مثال**: مقارنة قوائم المشتريات، الجينات، الكلمات المفتاحية •

6. II Hamming Distance

- للبيانات الفئوية والثنائية •
- المعادلة: عدد المواضع المختلفة •
- مسافة 1 = ("ABCD") و (ABCD") : مثال
- الاستخدام: البيانات الجينية، الرموز، البيانات الفئوية •

7. 📊 Mahalanobis Distance

- المعادلة: $(d = \sqrt{((x-\mu)^T \times \Sigma^{-1} \times (x-\mu))})$
- محسن يأخذ correlation محسن المتغيرات correlation
- القوة: يتعامل مع العلاقات بين المتغيرات •
- (بطيء) covariance matrix المشكلة: يحتاج حساب

8. Minkowski Distance

- المعادلة: $(d = (\Sigma | x_i y_i|^p)^{(1/p)})$
- :هو تعميم لـ
 - (p = 1): Manhattan
 - (p = 2): Euclidean
 - $(p = \infty)$: Chebyshev
- مرونة في التحكم بنوع المسافة •

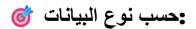
9. **Correlation Distance**

- (1 correlation_coefficient)
- المعنى: يقيس العلاقة الخطية، يتجاهل المقياس والمتوسط
- الاستخدام: البيانات الزمنية، تحليل الأنماط •

10. See Haversine Distance

- للمواقع الجغرافية •
- يحسب المسافة على سطح كروي (الأرض) •
- الاستخدام: GPS: الاستخدام

متی نستخدم کل مقیاس؟



بيانات مستمرة عادية:

- Euclidean (الأكثر شيوعاً)
- Manhattan (مع outliers)

بيانات عالية الأبعاد:

- **Cosine** (الأفضل)
- Manhattan
- تجنب Euclidean

:(TF-IDF) بيانات نصية

- Cosine (الأمثل)
- Euclidean (مقبول بعد normalization)

بيانات ثنائية (0/1):

- Jaccard (الأفضل)
- Hamming

بيانات فئوية:

- Hamming (الأمثل)
- Jaccard (للمجموعات)

:مواقع جغرافية

• Haversine (الوحيد المناسب)

حسب الخوارزمية 🔄

K-Means:

• افتراضي: **Euclidean**

• بديل: Manhattan, Cosine

DBSCAN:

يدعم: أي مقياس

• شائع: Euclidean, Manhattan

Hierarchical:

يدعم: جميع المقاييس •

• شائع: Euclidean, Cosine

KNN:

يدعم: أي مقياس

اختيار حسب نوع البيانات •

دسب خصائص البيانات

خاصية البيانات	الأفضل	تجنب	
Outliers کثیرة	Manhattan, Chebyshev	Euclidean	
أبعاد عالية	Cosine	Euclidean	
متغيرات مترابطة	Mahalanobis	Euclidean	
وحدات مختلفة	ابعد Scaling	scaling بدون	
بيانات متناثرة	Cosine, Jaccard	Euclidean	

(Curse of Dimensionality) تأثير الأبعاد العالية



افي الأبعاد العالية، جميع النقاط تصبح **متساوية المسافة** تقريباً

:مثال توضيحي

python

```
import numpy as np
from scipy.spatial.distance import pdist

# عند مختلفة في أبعاد مختلفة

for dims in [2, 10, 100, 1000]:

   data = np.random.normal(0, 1, (100, dims))
   distances = pdist(data, 'euclidean')
   print(f':باین {dims}: متوسط (distances.mean():.2f}, بیاد" (distances.var():.2f}")
```

والحل

- استخدم Cosine Distance
- أولاً PCA قلل الأبعادي. •
- كبديل Manhattan استخدم

أمثلة عملية

(سلوك العملاء) E-commerce (سلوك العملاء)

```
python

# [بيانات العملاء: [العمر, الدخل, عدد الطلبات, المبلغ الإجمالي]

customers = [[25, 50000, 12, 2400],
[30, 75000, 8, 3200],
[45, 100000, 15, 6000]]

# كانفل مقياس # أفضل مقياس أنفضل مقياس #
```

مثال 2: تحليل النصوص 📰

مثال 3: بيانات جينية (ثنائية) 🧇

نصائح للاختيار الصحيح

خطوات التحليل 🔍

- اعرف نوع بياناتك (مستمرة، فئوية، ثنائية) .1
- (outliers ،عادي، منحرف) تحقق من التوزيع .2
- اعرف عدد الأبعاد (قليل < 10، عالي > 50) 3.
- **جرب عدة مقاييس** وقارن النتائج .4

اختيار سريع 🔸

```
ایات عادیة → Euclidean

الهای outliers → Manhattan

ابعاد عالیة → Cosine

ابعاد عالیة → Cosine

ابعاد عالیة → Haversine
```

كود للاختبار 🧽

python

```
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
جرب مقاييس مختلفة #
distances = ['euclidean', 'manhattan', 'cosine']
results = {}
for distance in distances:
  if distance == 'cosine':
     Spectral أو Spherical K-Means استخدم #
    from sklearn.cluster import SpectralClustering
    model = SpectralClustering(n_clusters=3, affinity='cosine')
  else:
    model = KMeans(n_clusters=3, metric=distance)
  labels = model.fit_predict(X)
  score = silhouette_score(X, labels)
  results[distance] = score
print("أفضل مقياس:", max(results.items(), key=lambda x: x[1]))
```

التحسينات العملية

(2M records): للبيانات الكبيرة

- Euclidean/Manhattan: سریعان جداً
- **Cosine**: سریع مع sparse data
- Mahalanobis: (بطيء جداً)

وتوفير الذاكرة 💾

```
python

# بدلاً من float32 استخدم # بدلاً من float32 استخدم

X = X.astype(np.float32)

# للبيانات المتناثرة

from scipy.sparse import csr_matrix

X_sparse = csr_matrix(X)
```

نصائح للإنتاج 🍪

للبيانات الجديدة scaler احفظ معاملات الـ

- للبيانات الضخمة approximation استخدم
- **قس الوقت** واختر أسرع مقياس مناسب

outliers! 🎯 الخلاصة: Euclidean اللأبعاد العالية Manhattan مع