

(Outlier Detection) مقارنة خوارزميات اكتشاف القيم الشاذة

جدول المقارنة الرئيسي

الخاصية	LocalOutlierFactor	IsolationForest	EllipticEnvelope	StandardScaler
النوع	Density-based	Tree-based	Statistical	Preprocessing
الهدف	اكتشاف القيم الشاذة	اكتشاف القيم الشاذة	اكتشاف القيم الشاذة	تطبيع البيانات
المبدأ	Local density comparison	Random isolation	Gaussian distribution	Standardization
الافتراضات	توزيع محلي للبيانات	لا يحتاج افتراضات	توزيع طبيعي	-
التعقيد الحسابي	$O(n^2)$	$O(n \log n)$	$O(n^3)$	$O(n)$
أداء مع البيانات الكبيرة	بطيء	سريع	متوسط	سريع جداً
حساسية للمعايرة	نعم	لا	نعم	غير مطلوب

التفاصيل الفنية

1. LocalOutlierFactor (LOF)

- **كيف يعمل:** يقارن الكثافة المحلية لكل نقطة مع الكثافة المحلية لجيرانها
- (أعلى من 1 = قيمة شاذة) LOF Score: النتيجة
- المعاملات المهمة:
 - `n_neighbors`: عدد الجيران (افتراضي: 20)
 - `contamination`: نسبة القيم الشاذة المتوقعة

2. IsolationForest

- **كيف يعمل:** يبني أشجار عشوائية ويعزل النقاط الشاذة بسرعة
- **النتيجة:** Anomaly Score (قيم سالبة = قيم شاذة)
- **المعاملات المهمة:**
 - `n_estimators`: عدد الأشجار (افتراضي: 100)
 - `contamination`: نسبة القيم الشاذة المتوقعة
 - `max_samples`: عدد العينات لكل شجرة

3. EllipticEnvelope

- **كيف يعمل:** يفترض أن البيانات تتبع التوزيع الطبيعي المتعدد المتغيرات
- **النتيجة:** Binary classification (0/1)
- **المعاملات المهمة:**
 - `contamination`: نسبة القيم الشاذة المتوقعة
 - `support_fraction`: نسبة النقاط المستخدمة لحساب المركز

4. StandardScaler

- **الهدف:** تطبيع البيانات (mean=0, std=1)
- **ليس خوارزمية اكتشاف:** يُستخدم كخطوة تحضيرية
- **المعادلة:** $(x - \text{mean}) / \text{std}$

متى نستخدم كل خوارزمية؟

LocalOutlierFactor

استخدم عندما:

- البيانات لها كثافة متغيرة في مناطق مختلفة
- (local anomalies) تريد اكتشاف القيم الشاذة المحلية

- حجم البيانات صغير إلى متوسط ($> 10,000$ نقطة)
- الدقة أهم من السرعة

أمثلة:

- اكتشاف الاحتيال في المعاملات المصرفية
- تحليل شبكات التواصل الاجتماعي
- مراقبة جودة التصنيع

IsolationForest

استخدم عندما:

- حجم البيانات كبير جداً
- تريد سرعة في التنفيذ
- لا تعرف توزيع البيانات مسبقاً
- البيانات عالية الأبعاد

أمثلة:

- مراقبة الشبكات الحاسوبية
- اكتشاف البريد المزعج
- (Bioinformatics) تحليل البيانات الحيوية
- IoT مراقبة أنظمة

EllipticEnvelope

استخدم عندما:

- البيانات تتبع التوزيع الطبيعي تقريباً
- عدد الأبعاد محدود (> 50 بُعد)
- تريد نموذجاً إحصائياً قوياً

- البيانات نظيفة وخالية من الضوضاء

أمثلة:

- تحليل البيانات المالية (أسعار الأسهم)
- مراقبة العمليات الصناعية
- تحليل البيانات الطبية المعملية

StandardScaler

استخدم دائماً عندما:

- المتغيرات لها وحدات قياس مختلفة
- هناك فروق كبيرة في النطاقات
- قبل تطبيق خوارزميات تعتمد على المسافة

استراتيجية الاستخدام المتكاملة

مقترح Pipeline:

خوارزمية الاكتشاف → النتائج → StandardScaler → البيانات الخام

حسب حجم البيانات:

- EllipticEnvelope أو LOF ($< 1K$): صغير
- LOF أو IsolationForest ($1K-100K$): متوسط
- IsolationForest فقط ($> 100K$): كبير

حسب نوع البيانات:

- EllipticEnvelope: توزيع طبيعي
- IsolationForest: توزيع غير معروف
- LOF: بيانات معقدة التركيب

- أولاً StandardScaler: بيانات متنوعة الوحدات

نصائح عملية

⚠ تحذيرات مهمة:

1. **LOF**: بطيء مع البيانات الكبيرة، حساس لاختيار `n_neighbors`
2. **IsolationForest**: قد يفشل مع البيانات عالية الأبعاد جداً
3. **EllipticEnvelope**: يفترض التوزيع الطبيعي، حساس للقيم الشاذة في التدريب
4. **StandardScaler**: لا يزيل القيم الشاذة، فقط يطبع البيانات

💡 أفضل الممارسات:

- اختبر عدة خوارزميات وقارن النتائج
- بناءً على معرفتك بالمجال `contamination` استخدم
- قم بتصوير النتائج للتحقق من المعقولية
- اجمع بين عدة خوارزميات للحصول على نتائج أكثر قوة

:مثال عملي للاختيار

السيناريو: بيانات مبيعات شركة (100,000 معاملة، 15 متغير) الخيار الأمثل

1. StandardScaler لتطبيع البيانات
2. IsolationForest لاكتشاف المعاملات الشاذة (سرعة + دقة)