

Scalers & Preprocessing

مقارنة تقنيات المعايرة والتطبيع

جدول المقارنة الرئيسي

| الخاصية | StandardScaler | MinMaxScaler | RobustScaler | MaxAbsScaler | Normalizer |
|-----------------------------|----------------|--------------|---------------|--------------|------------|
| النطاق النهائي | mean=0, std=1 | [0, 1] | median=0, IQR | [-1, 1] | |
| المقاومة للـ Outliers | ضعيفة جداً | ضعيفة جداً | ممتازة | ضعيفة | متوسطة |
| يحافظ على التوزيع | نعم | نعم | نعم | نعم | لا |
| حساسية للقيم الشاذة | عالية جداً | عالية جداً | منخفضة جداً | عالية | متوسطة |
| سرعة التطبيق | سريع جداً | سريع جداً | سريع جداً | سريع جداً | سريع جداً |
| مناسب للتوزيع الطبيعي | ممتاز | جيد | ممتاز | جيد | جيد |
| مناسب للتوزيع المنحرف | ضعيف | ضعيف | جيد | ضعيف | متوسط |
| التعامل مع البيانات الجديدة | ممتاز | جيد | ممتاز | جيد | ممتاز |
| مناسب للخوارزميات | ممتاز | ممتاز | ممتاز | ممتاز | جيد |

| الخاصية | StandardScaler | MinMaxScaler | RobustScaler | MaxAbsScaler | Normalizer |
|--------------------------|----------------|--------------|--------------|--------------|------------|
| الخطية | | | | | |
| مناسب لخوارزميات المسافة | ممتاز | ممتاز | ممتاز | ممتاز | ممتاز |
| fit يحتاج على التدريب | نعم | نعم | نعم | نعم | لا |

التفاصيل الفنية

1. StandardScaler (Z-score normalization)

- المعادلة: $z = (x - \mu) / \sigma$
- النتيجة: متوسط = 0، انحراف معياري = 1
- الافتراضات: التوزيع الطبيعي مفضل
- الاستخدام:

```
python
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X_train)
```

2. MinMaxScaler

- المعادلة: $x_{scaled} = (x - min) / (max - min)$
- النتيجة: جميع القيم بين 0 و 1
- outliers المشكلة: حساس جداً للـ

- **الاستخدام:**

```
python

from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_scaled = scaler.fit_transform(X_train)
```

3. RobustScaler

- **المعادلة:** $x_scaled = (x - median) / IQR$
- **IQR:** Interquartile Range (Q75 - Q25)
- IQR و median لأنه يستخدم outliers **القوة:** مقاوم للـ
- **الاستخدام:**

```
python

from sklearn.preprocessing import RobustScaler
scaler = RobustScaler()
X_scaled = scaler.fit_transform(X_train)
```

4. MaxAbsScaler

- **المعادلة:** $x_scaled = x / |x_max|$
- **النتيجة:** القيم بين -1 و 1
- (القيم الصفرية) sparsity **المميز:** يحافظ على الـ
- (sparse data) **مناسب** لـ: البيانات المتناثرة

5. Normalizer

- **المعادلة:** $x_scaled = x / ||x||$
- **يعمل على:** كل صف منفرد (ليس العمود)
- **النتيجة:** magnitude = 1 كل صف له

- أنواع:
 - I_1 : 1 = المجموع المطلق
 - I_2 : 1 = الجذر التربيعي لمجموع المربعات
- مناسب لـ: معالجة النصوص، التشابه

6. QuantileTransformer

- أو طبيعي (uniform) المبدأ: يحول التوزيع إلى توزيع منتظم
- outliers القوة: مقاوم جداً للـ
- الخيارات:
 - `output_distribution='uniform'`: منتظم $[0,1]$
 - `output_distribution='normal'`: غاوسي
- بطيء نسبياً لكن فعال جداً

7. PowerTransformer

- الهدف: تحويل البيانات المنحرفة إلى توزيع طبيعي
- أنواع:
 - **Box-Cox**: للقيم الموجبة فقط
 - **Yeo-Johnson**: لأي قيم (موجبة وسالبة)
- مناسب لـ: البيانات شديدة الانحراف

8. Unit Vector Scaler (L2 Normalizer)

- **Normalizer** مع L2 norm نفس
- cosine similarity أو عند الحاجة لـ text data الاستخدام: عادة مع الـ

Scaler متى نستخدم كل

🎯 حسب نوع البيانات:

outliers: بيانات طبيعية التوزيع + بدون

- **StandardScaler** (الأفضل)

كثيرة outliers بيانات بها

- **RobustScaler** (الأمثل)
- **QuantileTransformer**

(skewed): بيانات منحرفة

- **PowerTransformer** (يصلح التوزيع)
- **QuantileTransformer**

(sparse): بيانات متناثرة

- **MaxAbsScaler**
- **MinMaxScaler/StandardScaler** تجنب

🤖 حسب نوع الخوارزمية:

خوارزميات تعتمد على المسافة:

- KNN, SVM, K-Means, Neural Networks
- **StandardScaler أو MinMaxScaler**: استخدم

خوارزميات خطية:

- Linear/Logistic Regression, Lasso, Ridge
- **StandardScaler**: استخدم

خوارزميات الأشجار:

- Random Forest, XGBoost, Decision Trees
- (لكن قد يحسن الأداء) **scaling** لا تحتاج

معالجة النصوص:

- استخدم: Normalizer (L2)

حسب شكل التوزيع

| شكل التوزيع | الأفضل | البديل |
|---------------|---------------------|---------------------|
| طبيعي | StandardScaler | MinMaxScaler |
| منتظم | MinMaxScaler | StandardScaler |
| منحرف | PowerTransformer | QuantileTransformer |
| outliers به | RobustScaler | QuantileTransformer |
| متعدد الأنماط | QuantileTransformer | RobustScaler |

أمثلة عملية

مثال 1: بيانات عادية

```
python

# بيانات الراتب والعمر - توزيع طبيعي تقريباً
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X_train)
```

outliers مثال 2: بيانات بها

```
python
```

بيانات الأسعار - بها قيم شاذة

```
from sklearn.preprocessing import RobustScaler
```

```
scaler = RobustScaler()
```

```
X_scaled = scaler.fit_transform(X_train)
```

مثال 3: بيانات منحرفة

python

بيانات الدخل - منحرفة جداً

```
from sklearn.preprocessing import PowerTransformer
```

```
scaler = PowerTransformer(method='yeo-johnson')
```

```
X_scaled = scaler.fit_transform(X_train)
```

مثال 4: بيانات نصية

python

TF-IDF vectors

```
from sklearn.preprocessing import Normalizer
```

```
scaler = Normalizer(norm='l2')
```

```
X_scaled = scaler.fit_transform(X_train)
```

نصائح مهمة للتطبيق

 أخطاء شائعة:

 خطأ شائع:

python

خطأ: scaling قبل train/test split

```
X_scaled = StandardScaler().fit_transform(X)
```

```
X_train, X_test = train_test_split(X_scaled)
```

✓ الطريقة الصحيحة:

python

```
# على التدريب فقط fit: صحيح
X_train, X_test = train_test_split(X)
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test) # فقط transform!
```

🎯 قواعد ذهبية:

1. على التدريب فقط **fit** دائماً
2. (histogram, boxplot) **حلل** التوزيع أولاً
3. scaling قبل الـ **outliers** تعامل مع الـ
4. وقارن النتائج **scalers** جرب عدة
5. لتجنب الأخطاء **Pipeline** استخدم

🔄 Pipeline مثالي:

python

```
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression

pipeline = Pipeline([
    ('scaler', StandardScaler()),
    ('model', LogisticRegression())
])

pipeline.fit(X_train, y_train)
predictions = pipeline.predict(X_test)
```


M Records تطبيق عملي لـ 2

(2M) للبيانات الضخمة:

✓ **السرعة Scalers** الـ:

1. **StandardScaler**: الأسرع والأكثر استقراراً
2. **MinMaxScaler**: سريع جداً
3. **RobustScaler**: outliers سريع ومقاوم للـ

⚠ **تجنب مع البيانات الكبيرة**:

- **QuantileTransformer**: بطيء جداً
- **PowerTransformer**: أبطأ نسبياً

🚀 **نصائح للأداء**:

```
python
```

```
# للبيانات الكبيرة
```

```
from sklearn.preprocessing import StandardScaler
```

```
import numpy as np
```

```
# لتوفير الذاكرة dtype=float32 استخدم
```

```
scaler = StandardScaler()
```

```
X_scaled = scaler.fit_transform(X.astype(np.float32))
```

اختيار سريع حسب الحالة

🎯 **الاختيار السريع**:

(آمن دائماً) **StandardScaler** → لا أعرف شيء عن البيانات

RobustScaler → واضحة outliers البيانات بها

MinMaxScaler → أريد كل شيء بين 0 و 1

PowerTransformer → البيانات منحرفة جداً

Normalizer → similarity: للنصوص أو

QuantileTransformer (لو الوقت يسمح) → أريد أقصى دقة ممكنة

🎯! هو الخيار الآمن لـ 90% من الحالات StandardScaler: الخلاصة