

在特征工程中，异常值（特别大或特别小的数值）确实可能导致归一化表现不佳。例如，在使用 **Min-Max Scaling** 或 **Z-Score Standardization** 时，异常值会极大地影响缩放范围或均值/标准差，从而导致模型性能下降。因此，处理异常值是特征工程中的重要一步。以下是一些常见的处理方法和具体建议：

1. 检测异常值的方法

在处理异常值之前，首先需要检测哪些是异常值：

1.1 基于统计规则

- **IQR 法（四分位数间距）：**

- 计算特征的 **第一四分位数 (Q1)** 和 **第三四分位数 (Q3)**。
- 定义异常值范围：低于 $Q1 - 1.5 \times IQR$ 或高于 $Q3 + 1.5 \times IQR$ 的数据为异常值， $IQR = Q3 - Q1$ 。
- 示例代码：

```
import numpy as np
Q1 = np.percentile(data, 25)
Q3 = np.percentile(data, 75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
outliers = data[(data < lower_bound) | (data > upper_bound)]
```

- **3 σ 法则（正态分布假设）：**

- 如果数据服从正态分布，取均值 μ 和标准差 σ ，定义异常值为低于 $\mu - 3\sigma$ 或高于 $\mu + 3\sigma$ 的数据。
- 示例代码：

```
mean = np.mean(data)
std = np.std(data)
lower_bound = mean - 3 * std
upper_bound = mean + 3 * std
outliers = data[(data < lower_bound) | (data > upper_bound)]
```

1.2 数据可视化

- **箱线图 (Box Plot)**：直观显示数据的分布及异常值。

```
import matplotlib.pyplot as plt
import seaborn as sns
sns.boxplot(data)
plt.show()
```

- **直方图 (Histogram)**：观察数据的整体分布，发现极端值。

```
plt.hist(data, bins=30)
plt.show()
```

- **散点图 (Scatter Plot)**：对于多维数据，可以绘制散点图查看是否存在异常点。

2. 处理异常值的方法

2.1 删除异常值

- **适用场景：**
 - 数据量较大，删除少量异常值不会导致数据量不足。
 - 异常值是由采集错误或录入错误导致，明确不属于真实数据。
- **方法：**
 1. 根据检测到的异常值范围，直接剔除这些数据。

2. 示例代码：

```
clean_data = data[(data >= lower_bound) & (data <= upper_bound)]
```

- **注意事项：**

- 如果异常值是重要信息（例如银行欺诈检测中的异常交易），则不能简单删除。
-

2.2 替换异常值

- **方法 1：用统计量替换**

- 用 **均值** 或 **中位数** 替换异常值。
- 示例代码：

```
data[data < lower_bound] = np.median(data)  
data[data > upper_bound] = np.median(data)
```

- **方法 2：用分位数替换**

- 将异常值替换为特征的 **最小分位数**（如 **Q1**）或 **最大分位数**（如 **Q3**）。
- 示例代码：

```
data[data < lower_bound] = Q1  
data[data > upper_bound] = Q3
```

- **方法 3：用模型预测值替换**

- 使用回归模型或插值方法预测异常值，然后用预测值替换。
- 示例代码（线性插值）：

```
data = data.interpolate(method='linear')
```

2.3 对异常值进行截断 (Winsorization)

- 原理：
 - 将异常值限制在一个合理范围内，而不是直接删除或替换。
 - 例如：将低于 $Q1 - 1.5 \times IQR$ 的值设置为 $Q1 - 1.5 \times IQR$ ，将高于 $Q3 + 1.5 \times IQR$ 的值设置为 $Q3 + 1.5 \times IQR$ 。

- 示例代码：

```
data = np.clip(data, lower_bound, upper_bound)
```

- 优点：
 - 保留了所有数据点，避免删除数据。
 - 对模型的训练影响较小。
-

2.4 使用鲁棒的归一化方法

- 如果异常值无法完全消除，可以使用对异常值不敏感的归一化方法。

(1) RobustScaler (Sklearn 提供)

- 通过使用 **中位数 (Median)** 和 **四分位数范围 (IQR)** 进行缩放，而不是均值和标准差。

- 公式：

$$x_{\text{scaled}} = \frac{x - \text{median}(x)}{Q3 - Q1}$$

- 示例代码：

```
from sklearn.preprocessing import RobustScaler  
scaler = RobustScaler()
```

```
scaled_data = scaler.fit_transform(data.reshape(-1, 1))
```

(2) Log 变换

- 对于分布偏斜或有较大异常值的数据，可以对其取对数。对数变换会压缩大值的范围。
- 示例代码：

```
import numpy as np
data = np.log1p(data) # log(1 + x), 避免 log(0) 问题
```

(3) Box-Cox 变换

- Box-Cox 变换是一种对数据进行正态化的变换方法，适用于正值数据。
- 示例代码：

```
from scipy.stats import boxcox
data_transformed, _ = boxcox(data)
```

(4) Yeo-Johnson 变换

- 类似于 Box-Cox，但适用于负值和零值数据。
- 示例代码：

```
from sklearn.preprocessing import PowerTransformer
transformer = PowerTransformer(method='yeo-johnson')
transformed_data = transformer.fit_transform(data.reshape(-1, 1))
```

2.5 特殊场景下的异常值处理

- 时间序列数据：

- 对于时间序列中的异常值，可以使用滑动窗口或季节性分解方法检测并替换异常值。
- 示例：

```
from statsmodels.tsa.seasonal import seasonal_decompose
result = seasonal_decompose(data, model='additive', period=12)
data_cleaned = data - result.resid # 去除残差
```

- 分类任务中的异常值：

- 如果异常值是某些类别的重要特征（如欺诈检测中的异常交易），应保留异常值，并通过特征工程将其显式建模。

3. 归一化后的验证

处理异常值后，再次归一化数据并验证归一化效果是否改进。常见的验证方法包括：

- 可视化数据分布（如绘制直方图或箱线图）。
- 观察归一化后的均值和标准差是否合理。
- 检查模型训练效果是否提升。

4. 总结

处理异常值的常用流程：

1. 检测异常值：

- 使用统计规则（如 IQR 或 3σ 法）或可视化方法（如箱线图、散点图）。

2. 选择合适的方法处理异常值：

- 数据量大且异常值明确无效 → 删除异常值。
- 异常值有效但对模型有负面影响 → 替换或截断异常值。
- 异常值无法有效处理 → 使用鲁棒的归一化方法（如 RobustScaler）。

3. 验证归一化效果：

- 通过可视化或统计指标验证归一化是否合理。

异常值的处理需要结合具体任务场景，避免简单粗暴地删除或替换，特别是在异常值可能包含有用信息的情况下。