① read data ② clean data ③ remove duplicated data
④ outliers (Detecting and Filtering) ✓
⑤ Feature selection ⑥ Data normalization ✓ □ = 1

4 ≠ outliers detecting    outliers الـ عن تشرح

                                    ✓ ورَ
                                    العقد ⟸

إننا بالتقى قيم خارج الـ Range وهي ضعيفة لبعض أسباب
أمن متوسط الأعمار 65 سنة وأتمكن عمرقد يصل إلى 100 عام
لكن ماذا عن وجدت أعمار تتجاوز الـ 500 والـ 700 هذا غير منطقي
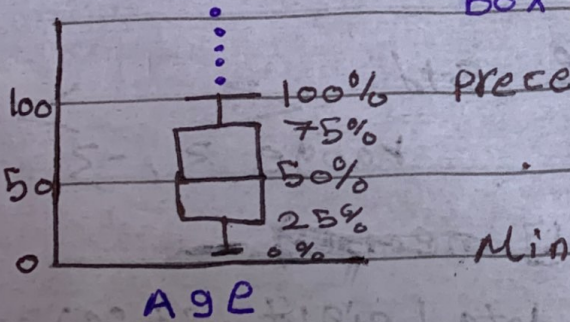فهذه القيمة تسمى outliers

التعرف على هذه القيم بعد طرق ⟸

1) عن طريق Domain knowledge
تفترض إنت الأعمار تصل إلى 100 فأي شي أكثر
age ≤ 105 إلى يتجاوز الرقم يُعتبر outliers

2) عن طريق الـ visualization
                Box Plot, Bar Plot
                precentile   Max

100 ─────── 100%
        75%
50 ─────── 50%
        25%
0 ──────── 0%    Min

    Age

Subject : _____          Date:_____

[1] Identify outliers by code          1- Domain
    1- Identify Max-Val, min-Val             knowledge

مش بنعتمد علي اي حاجة فوق او اخت فلام
[1] Max_Value = data['height'].quantile(0.95)
    8.37 inch  بي لو نحسب الارتفاع فوق قيمة معينة
[2] Min_Value = data['height'].quantile(0.05)
    2.3             أقل قوة يسيب حجم بي فوق في الطول

[3] Show The outliers [Identify outliers by domain]
    data[ data['height'] < min-value]
    data[ data['height'] > max_value]

| name  | highet |
|-------|--------|
| essra | 0.9    |

| name    | height |
|---------|--------|
| amr     | 14.9   |
| mostafa | 13.2   |
| Ibrahim | 15.3   |

[4] Remove outliers for simple data

df [ (df['height'] < max_value) & (df['height'] > min_vale)]
احنا حددنا دي هي القيم اللي مش عايزنها outliers

جوة جي جيم_جي
quntile