

Outlier

Identify and dealing with outliers is important in various field including statistics, ML, data analysis.

ସି. ଗ୍ରା (2.4) ରେ, $\mu_{\text{result}} = 1.5$ ଏବଂ $\sigma = 0.5$ ରହିବ।

theory.

• Measurement Errors: mistakes in record data on entry data
→ 1000 Enter થયેલ 100 વચ્ચે (પ્રભ)

DATA entry errors; Human or system errors while inputting data.
- Misplace decimal - numerical data.

Sampling error: Data collected from a non representative sample
gives down or give wrong answer.

Data processing issues: Errors in data transmission, cleaning
- Duplicate records, mismatched data.

Why outlier Important

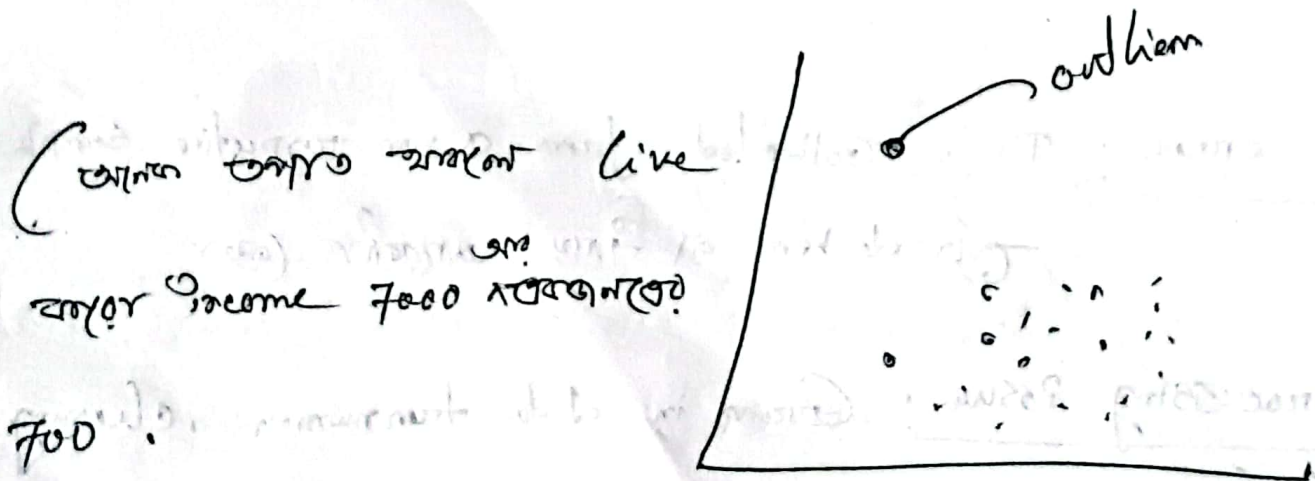
- Affective model performance, Algorithm like, Linear Regression, k-means, NN are especially sensitive to outliers
- Outliers heavily influence the mean, Standard deviation

Summaries.

- Outliers may reveal unexpected pattern
- Proper handling of outliers ensure cleaner data.

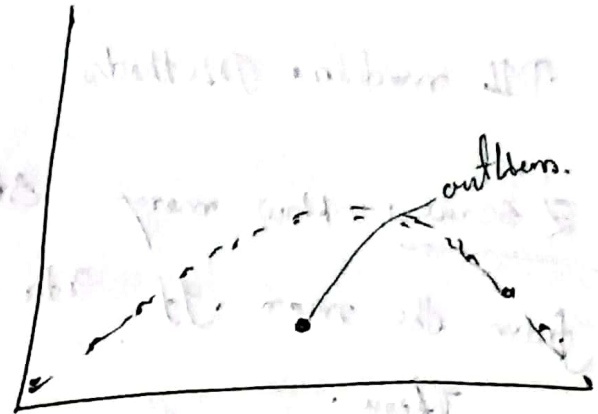
Types of outliers

Global outliers: It is a data points that are significantly different from all other data points in the dataset.
Outliers easily identify.



Contextual/Conditional outliers: It is data points that are considered only in specific context or condition.

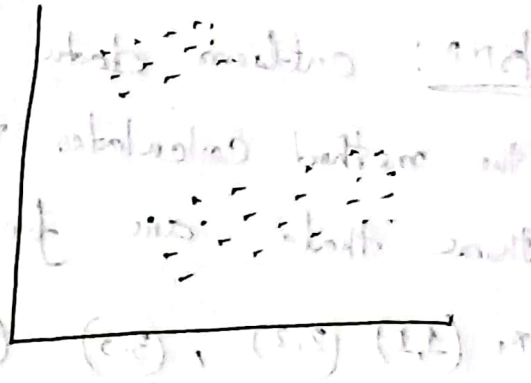
- Require additional context to identify
- often linked to time, loc, other conditions.



Collective outliers: It is a group of data points that together differ from the overall pattern, even individual points might not be outliers.

- indicate different behaviour on pattern

— Long distance to different group



Example CGPA

1st = 2.2 2nd = 2.3 3rd = 2.2

4th = 3.57 5th = 3.83 6th = 3.81



Outliers Detection

Outliers Detection ensured the quality and accuracy of ML models. Methods.

Z score: How many standard deviation a data point is from the mean. If it's greater than threshold (3, -3) considered as outlier.

Ex $\rightarrow (80, 85, 88, 90, 92, 100, 250)$ mean = 96.42.
Standard deviation = 58.79.

$$Z = \frac{250 - 96.42}{58.79} = 2.61 < 3, \text{ so } 250 \text{ is outlier.}$$

RNN: Outliers data points that have few neighbors. The method calculates the distance from each point, identify those that are further away than others.

Ex: (1,1) (2,2), (3,3) (100,100) \leftarrow (100,100) is outlier.

IQR: Interquartile Range

Calculate the range 1st quartile (Q_1) and 3rd quartile.

Ex: $\rightarrow [1, 2, 3, 4, 5, 6, 7, 100]$ $Q_1 = 2.5$, $Q_3 = 6.5$

$$IQR = 6.5 - 2.5 = 4, \quad 1.5 \times 4 = 1.5$$

Upper bound = $6.5 + 1.5 \times 4 = 13.5 < 100$ So 100 is outlier

Local outlier factor (LOF): Measure density of data points.

- If a point has lower density than its neighbors, it is an outlier.

(1,1), (2,2), (3,3) (50,50) ← outlier

Clustering → DBSCAN: It groups points that are close to each other and identifies outliers as points that don't belong to any other cluster.

(1,1), (2,2), (10,50), (100,100) → high density
→ outlier because it is far from the center of the cluster.

k-nearest cluster: Points that are far from the center of their cluster are considered outliers.

Outliers Detection

Supervised

- Z-Score
- modified Z-Score
- Inter Quartile Range
- SVM
- ANN

Unsupervised

- Elliptic Envelope
- Isolation forest
- k-means
- DBSCAN
- Local outlier factor