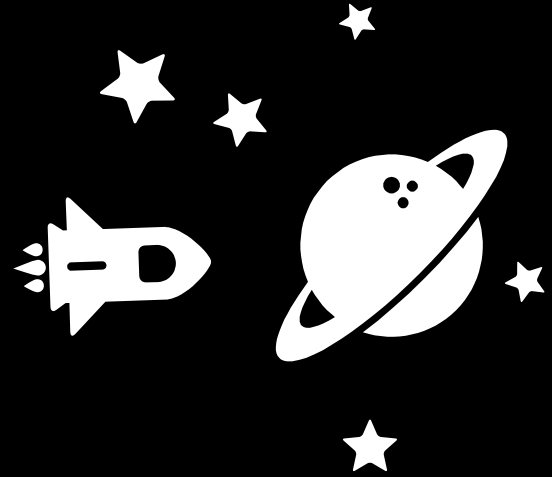


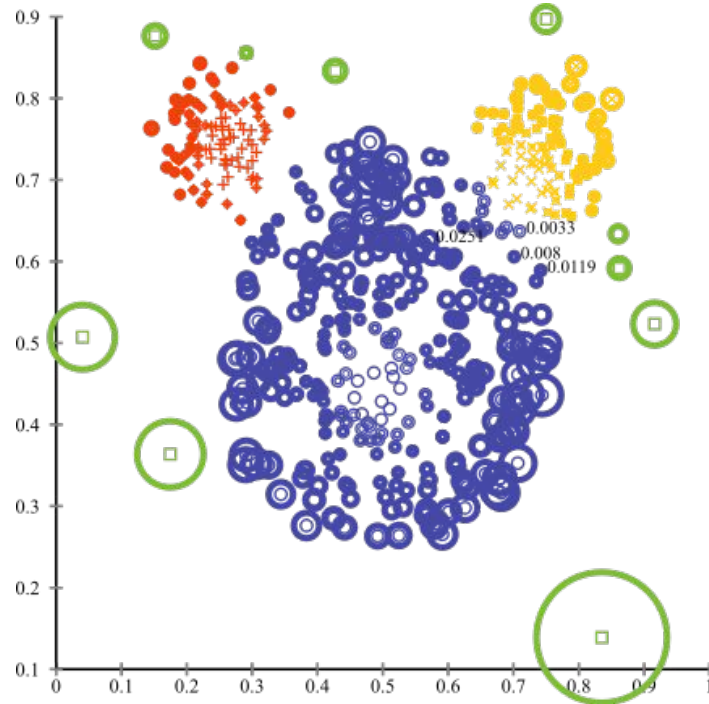
# **Outliers: Analysis & Detection**

**Presentation by Berk Sudan**

# Examples

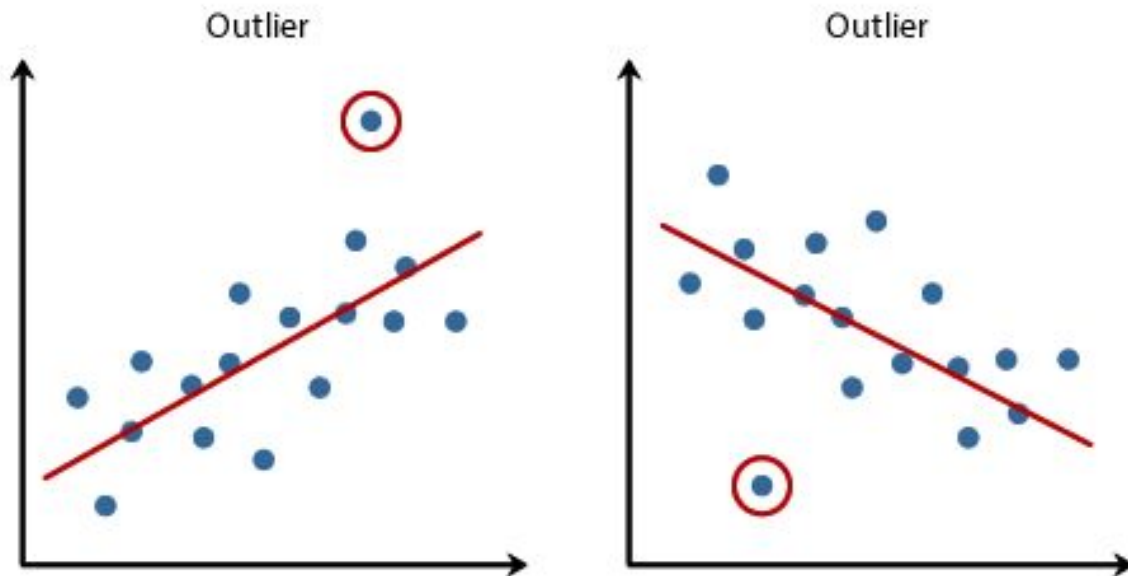


# Clustering



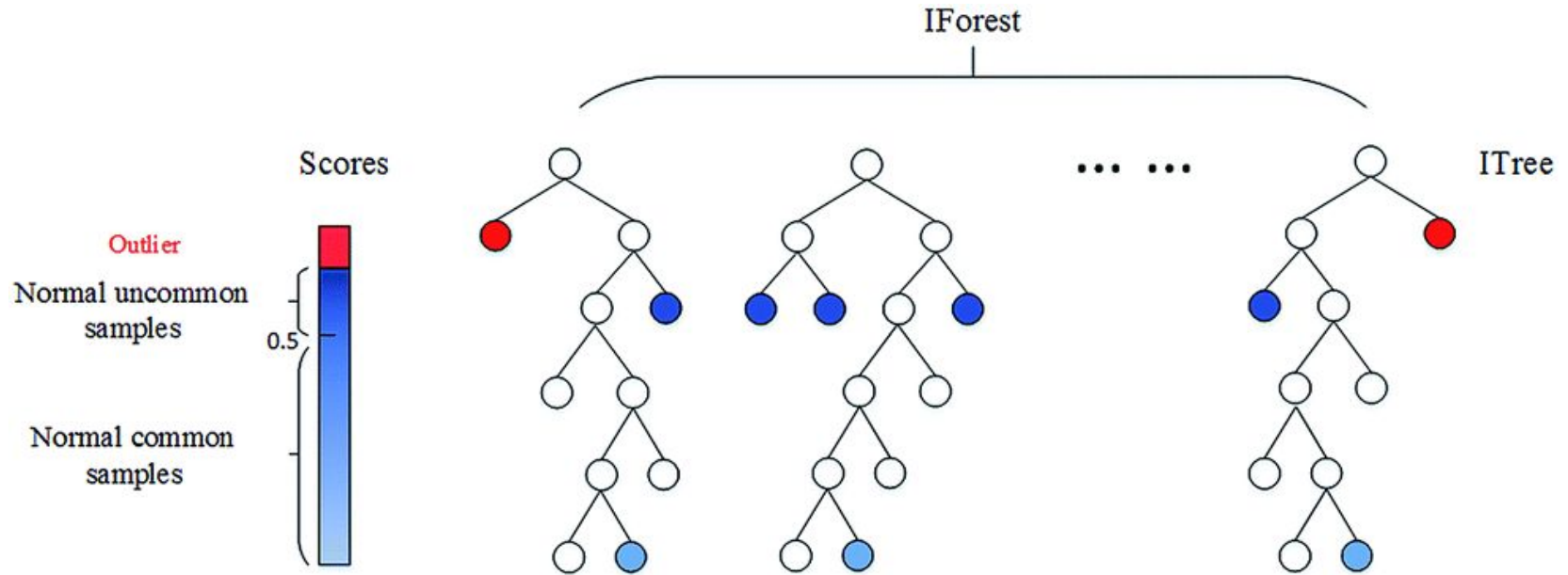
**Ref:** <https://stats.stackexchange.com/questions/160260/anomaly-detection-based-on-clustering>

# Regression



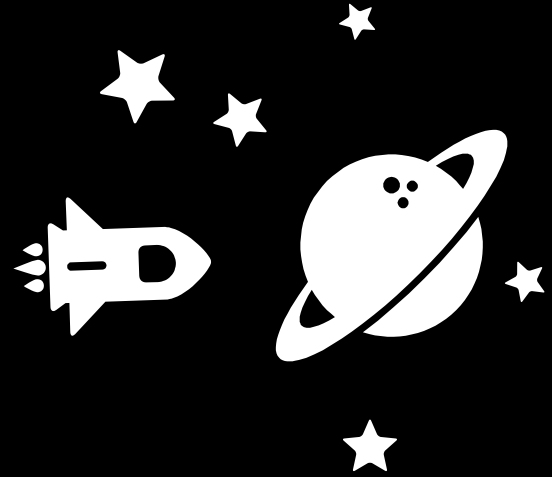
*Ref: <https://datanee.com/2016/08/11/outlier-detection-an-overview-and-applications/>*

# Trees



Ref: <https://towardsdatascience.com/a-brief-overview-of-outlier-detection-techniques-1e0b2c19e561>

**What can  
outliers  
possibly be?**



## Possibility of Invalid Data - 1

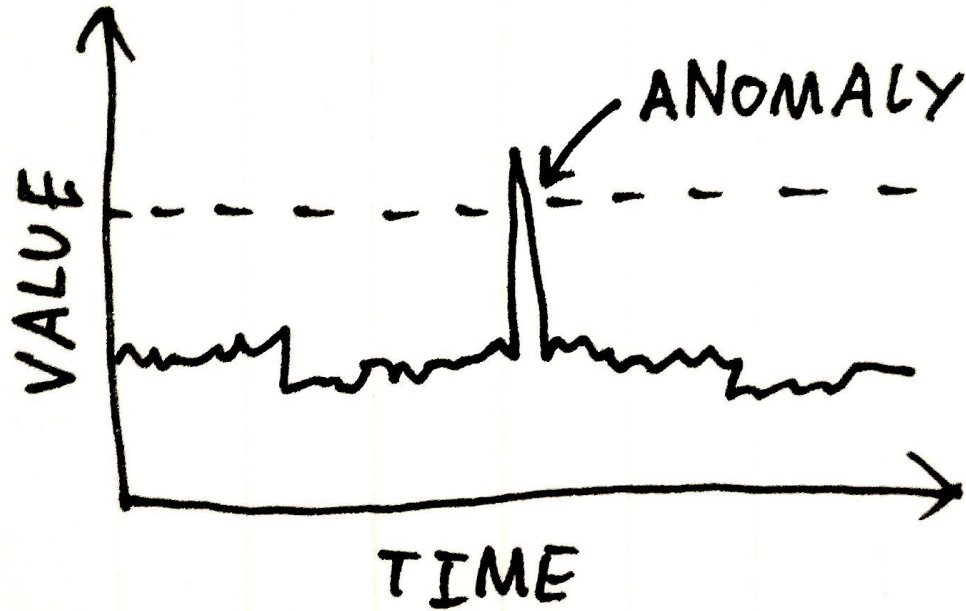
Name	Gender	Year of Birth
-----	-----	-----
Andrew NG	M	1976
Sarah Tan	F	2976

## Possibility of Invalid Data - 2

Gender	isPregnant	Age	isSmoking?
Male	Yes	25	Yes
Female	Yes	25	Yes
Male	No	25	No

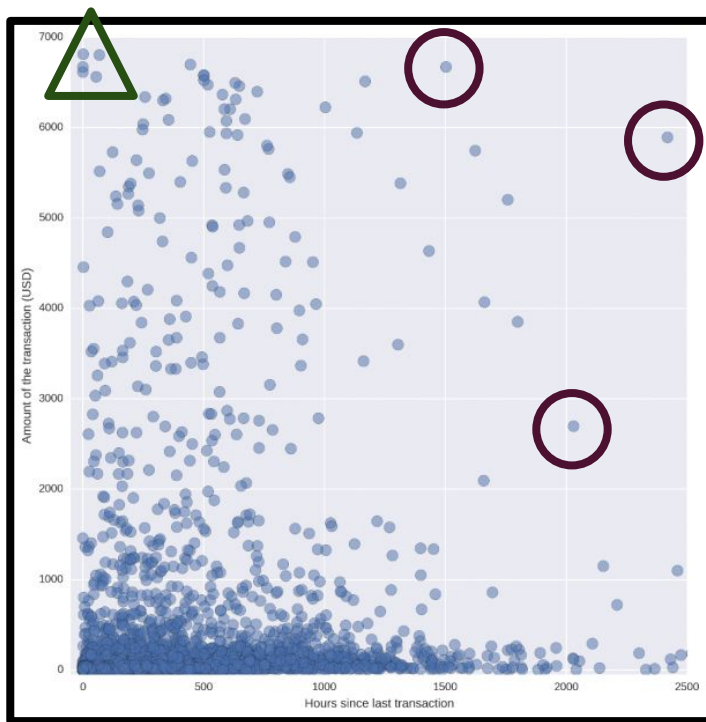


# Anomaly



Ref: <https://medium.com/the-data-dynasty/anomaly-detection-in-google-analytics-a-new-kind-of-alerting-9c31c13e5237>

# Fraud



**X axis:** Hours since last transaction

**Y axis:** Amount of the transaction (in USD)

○ : Outliers but not fraud

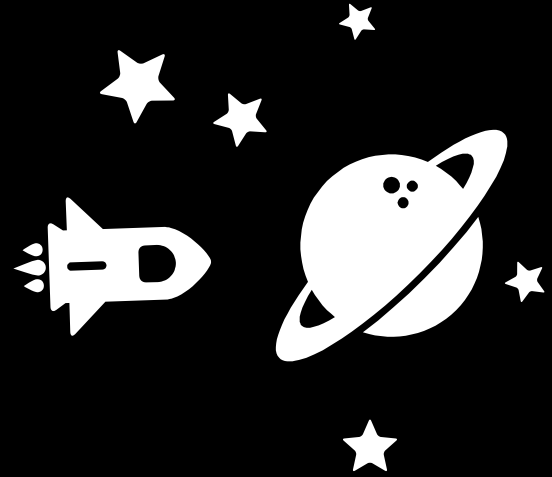
△ : Outliers & fraud

*Ref: <https://blog.easysol.net/advanced-outlier-detection/>*

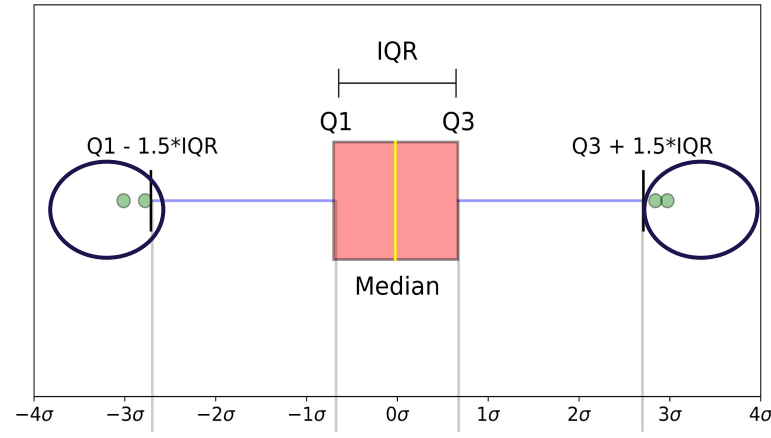
## What else?

- > **Big variance, shifted mean** due to outliers
- > **Noise** in data
- > **Highly skewed** data distribution
- > **Additional problems** in Distance-Based Algorithms  
(e.g, K-Means, K-NN, ...)

# Methods: Uni-variate



# Using Box-Plot



**For value  $X=x$ :**

$(x > (Q3 + 1.5 \cdot IQR)) \vee (x < (Q1 - 1.5 \cdot IQR)) \rightarrow (x \text{ is an outlier})$

$(x > (Q3 + 3 \cdot IQR)) \vee (x < (Q1 - 3 \cdot IQR)) \rightarrow (x \text{ is an extreme-value})$

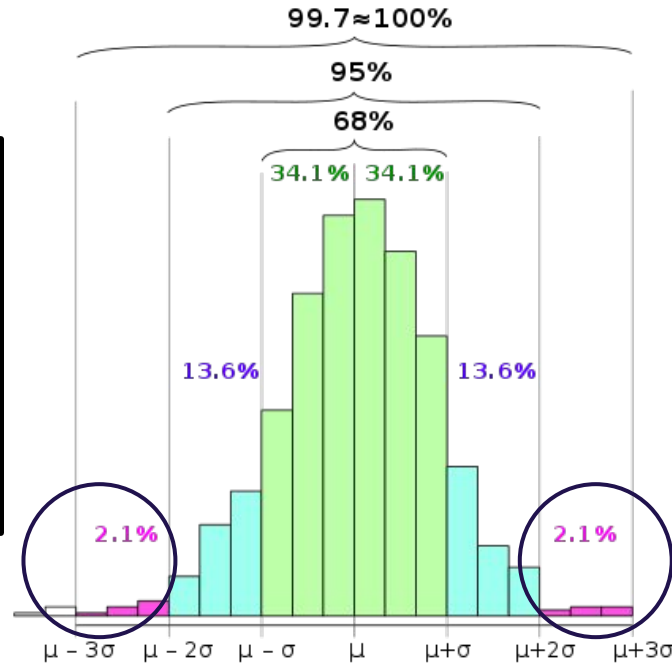
*Ref: <https://towardsdatascience.com/understanding-boxplots-5e2df7bcbd51>*

# Using Standard Deviation

For value  $X=x$ :

$(x > (\mu + 2\sigma)) \vee (x < (\mu - 2\sigma)) \rightarrow$   
(x is an **outlier**)

$(x > (\mu + 3\sigma)) \vee (x < (\mu - 3\sigma)) \rightarrow$   
(x is an **extreme-value**)



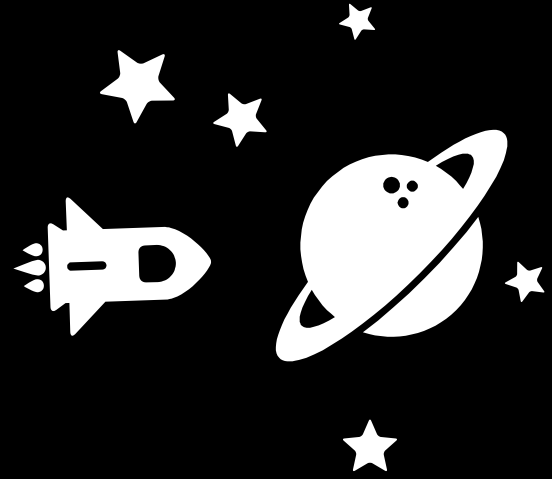
○ : outliers

Ref: [https://en.wikipedia.org/wiki/68%E2%80%939395%E2%80%939399.7\\_rule](https://en.wikipedia.org/wiki/68%E2%80%939395%E2%80%939399.7_rule)

## Hard-Edges Method:

Data yielding outside of the (***1<sup>th</sup>* - 99<sup>th</sup>**)  
**quantile/percentile interval** will be  
evaluated **as *outlier***.

# Methods: Multi-variate





# Using Distance based Approach

Euclidean Distance:

$$d(p, q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + (p_3 - q_3)^2}$$

**FOR** any point:

**IF:** distance (d)  $\geq$  threshold\_distance **AND**

fraction for non-neighbour sample (f)  $\geq$  threshold\_fraction

**THEN: outlier!!**

# Using Distance based Approach - Example

$$P_1 = (2,4)$$

$$P_2 = (3,2)$$

$$P_3 = (1,1)$$

$$P_4 = (4,3)$$

$$P_5 = (1,6)$$

$$P_6 = (5,3)$$

$$P_7 = (4,2)$$

Find the outliers using the distance-based technique.

- The threshold distance is 3 (i.e,  $d \geq 3$ )
- The threshold fraction  $f$  for non-neighbour sample is  $4/7$  (i.e,  $f \geq 4/6$ )

# Using Distance-based Approach - Example

Distance Matrix for Euclidean Distance

$$D_{7 \times 7} = \begin{array}{c} \begin{array}{ccccccc} p_1 & p_2 & p_3 & p_4 & p_5 & p_6 & p_7 \end{array} \\ \left[ \begin{array}{ccccccc} 0 & \sqrt{5} & \sqrt{10} & \sqrt{5} & \sqrt{5} & \sqrt{10} & \sqrt{8} \\ . & 0 & \sqrt{5} & \sqrt{2} & \sqrt{20} & \sqrt{5} & 1 \\ . & . & 0 & \sqrt{13} & 5 & \sqrt{20} & \sqrt{10} \\ . & . & . & 0 & \sqrt{18} & 1 & 1 \\ . & . & . & . & 0 & 5 & 5 \\ . & . & . & . & . & 0 & \sqrt{2} \\ . & . & . & . & . & . & 0 \end{array} \right] \begin{array}{c} p_1 \\ p_2 \\ p_3 \\ p_4 \\ p_5 \\ p_6 \\ p_6 \end{array} \end{array}$$

# Using Distance-based Approach - Example

Distance Matrix for Euclidean Distance

- The threshold distance is 3 (i.e,  $d \geq 3$ )
- The threshold fraction  $f$  is  $4/6$  (i.e,  $f \geq 4/6$ )

$$D_{7 \times 7} = \begin{matrix} & \begin{matrix} p_1 & p_2 & p_3 & p_4 & p_5 & p_6 & p_7 \end{matrix} \\ \begin{matrix} p_1 \\ p_2 \\ p_3 \\ p_4 \\ p_5 \\ p_6 \\ p_7 \end{matrix} & \begin{bmatrix} 0 & \sqrt{5} & \sqrt{10} & \sqrt{5} & \sqrt{5} & \sqrt{10} & \sqrt{8} \\ . & 0 & \sqrt{5} & \sqrt{2} & \sqrt{20} & \sqrt{5} & 1 \\ . & . & 0 & \sqrt{13} & 5 & \sqrt{20} & \sqrt{10} \\ . & . & . & 0 & \sqrt{18} & 1 & 1 \\ . & . & . & . & 0 & 5 & 5 \\ . & . & . & . & . & 0 & \sqrt{2} \\ . & . & . & . & . & . & 0 \end{bmatrix} \end{matrix} \begin{matrix} p_1 \rightarrow f = 2 \\ p_2 \rightarrow f = 1 \\ p_3 \rightarrow f = 5 \\ p_4 \rightarrow f = 2 \\ p_5 \rightarrow f = 5 \\ p_6 \rightarrow f = 3 \\ p_7 \rightarrow f = 2 \end{matrix}$$

# Using Distance-based Approach - Example

Distance Matrix for Euclidean Distance

- The threshold distance is 3 (i.e,  $d \geq 3$ )
- The threshold fraction  $f$  is  $4/6$  (i.e,  $f \geq 4/6$ )

$$D_{7 \times 7} = \begin{array}{ccccccc|l} p_1 & p_2 & p_3 & p_4 & p_5 & p_6 & p_7 & \\ \hline 0 & \sqrt{5} & \sqrt{10} & \sqrt{5} & \sqrt{5} & \sqrt{10} & \sqrt{8} & p_1 \rightarrow f = 2 \\ \cdot & 0 & \sqrt{5} & \sqrt{2} & \sqrt{20} & \sqrt{5} & 1 & p_2 \rightarrow f = 1 \\ \cdot & \cdot & 0 & \sqrt{13} & 5 & \sqrt{20} & \sqrt{10} & p_3 \rightarrow f = 5 \\ \cdot & \cdot & \cdot & 0 & \sqrt{18} & 1 & 1 & p_4 \rightarrow f = 2 \\ \cdot & \cdot & \cdot & \cdot & 0 & 5 & 5 & p_5 \rightarrow f = 5 \\ \cdot & \cdot & \cdot & \cdot & \cdot & 0 & \sqrt{2} & p_6 \rightarrow f = 3 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & 0 & p_7 \rightarrow f = 2 \end{array}$$

So, the outliers are:  $\{P_3, P_5\}$

# **End of Presentation**

**Presented by Berk Sudan**