

---

# DATA ANALYSIS

OUTLIERS AND ANOMALOUS OBSERVATIONS

# ANOMALOUS OBSERVATIONS

In practice, an **anomalous observation** may arise as

- a **“bad” object/measurement**: data artifacts, spelling mistakes, poorly imputed values, etc.
- a **misclassified observation**: according to the existing data patterns, the observation should have been labeled differently;
- an observation whose measurements are found in the **distribution tails** of a large enough number of features;
- an **unknown unknown**: a completely new type of observations whose existence was heretofore unsuspected.

---

# ANOMALOUS OBSERVATIONS

---

Observations could be anomalous in one context, but not in another:

- a 6-foot tall adult male is in the 86th percentile for **Canadian males** (tall, but not unusual);
- in **Bolivia**, the same man would be in the 99.9th percentile (very tall and unusual).

Anomaly detection points towards interesting questions for analysts and SMEs: in this case, **why is there such a large discrepancy** in the two populations?

# OUTLIERS

**Outlying observations** are data points which are **atypical** in comparison to

- the unit's remaining features (*within-unit*),
- the field measurements for other units (*between-units*)

Outliers are observations which are **dissimilar to other cases** or which **contradict known dependencies** or rules.

Careful study is needed to determine whether outliers should be retained or removed from the dataset.

# DETECTING ANOMALIES

Outliers may be anomalous along any of the unit's variables, or in combination.

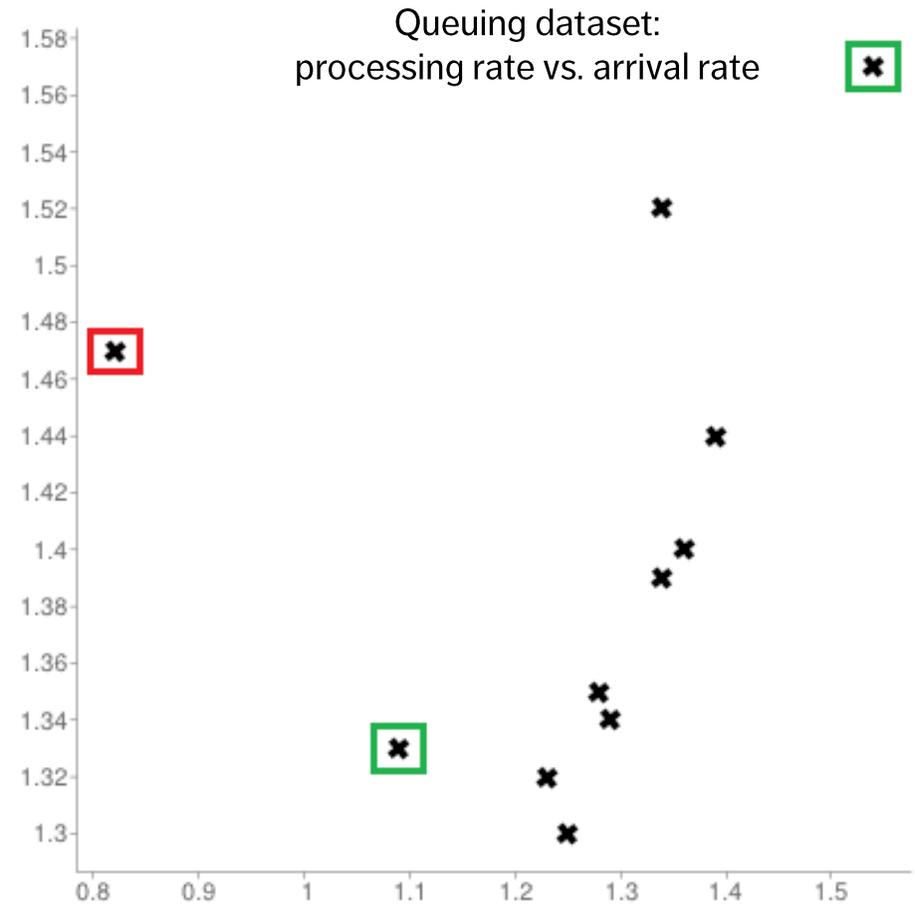
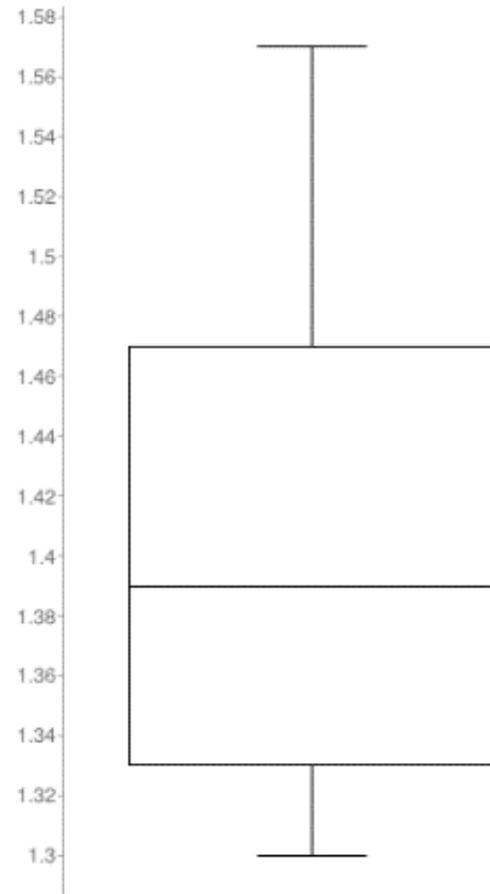
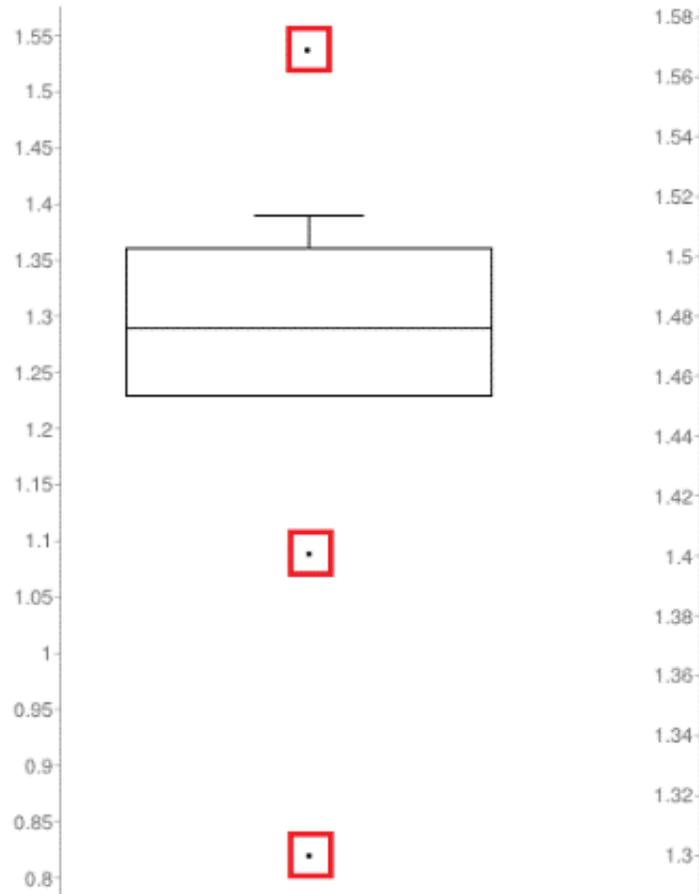
Anomalies are by definition **infrequent** and typically shrouded in **uncertainty** due to small sample sizes.

Differentiating anomalies from noise or data entry errors is **hard**.

Boundaries between normal and deviating units may be **fuzzy**.

Anomalies associated with malicious activities are typically **disguised**.

# VISUAL OUTLIER DETECTION



# DETECTING ANOMALIES

Numerous methods exist to identify anomalous observations; **none of them are foolproof** and judgement must be used.

Graphical methods are easy to implement and interpret:

- **outlying observations**

box-plots, scatterplots, scatterplot matrices, 2D tour, Cooke's distance, normal qq plots

- **influential data**

some level of analysis must be performed (leverage)

**Careful:** once anomalous observations have been removed from the dataset, previously “regular” units may become anomalous.

# ANOMALY DETECTION ALGORITHMS

**Supervised methods** use a historical record of labeled anomalous observations:

- domain expertise is required to tag the data
- classification or regression task
- rare occurrence problem

		Predicted Class	
		Normal	Anomaly
Actual Class	Normal	<i>TN</i>	<i>FP</i>
	Anomaly	<i>FN</i>	<i>TP</i>

**Unsupervised methods** don't use external information:

- traditional methods and tests
- can also be seen as a clustering or association rules problem

# ANOMALY DETECTION ALGORITHMS

The mis-classification cost is often assumed to be symmetrical, which can lead to **technically correct but useless** outputs.

For instance, most (99.999+%) air passengers do not bring weapons with them on flights; a model that predicts that no passenger is smuggling a weapon would be 99.999+% accurate, but it would miss the point completely.

For the **security agency**, the cost of wrongly thinking that a passenger is:

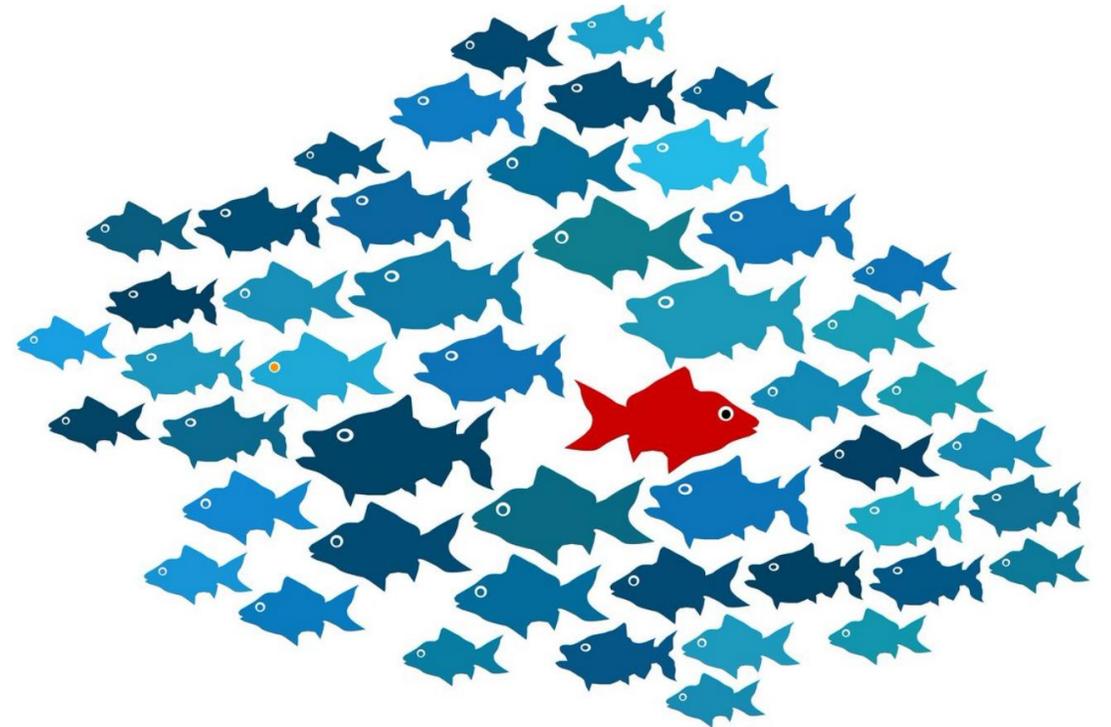
- smuggling a weapon  $\Rightarrow$  cost of a single search
- NOT smuggling a weapon  $\Rightarrow$  catastrophe (potentially)

But **wrongly targeted individuals** may have a different take on this!

# ANOMALY DETECTION ALGORITHMS

If all participants in a workshop except for one can view the video conference lectures, then the one individual/internet connection/computer is **anomalous** – it behaves in a manner which is different from the others.

But this **DOES NOT MEAN** that the different behaviour is necessarily the one we are interested in...



# INFLUENTIAL OBSERVATIONS



**Influential data points** are observations whose absence leads to **markedly different** analysis results.

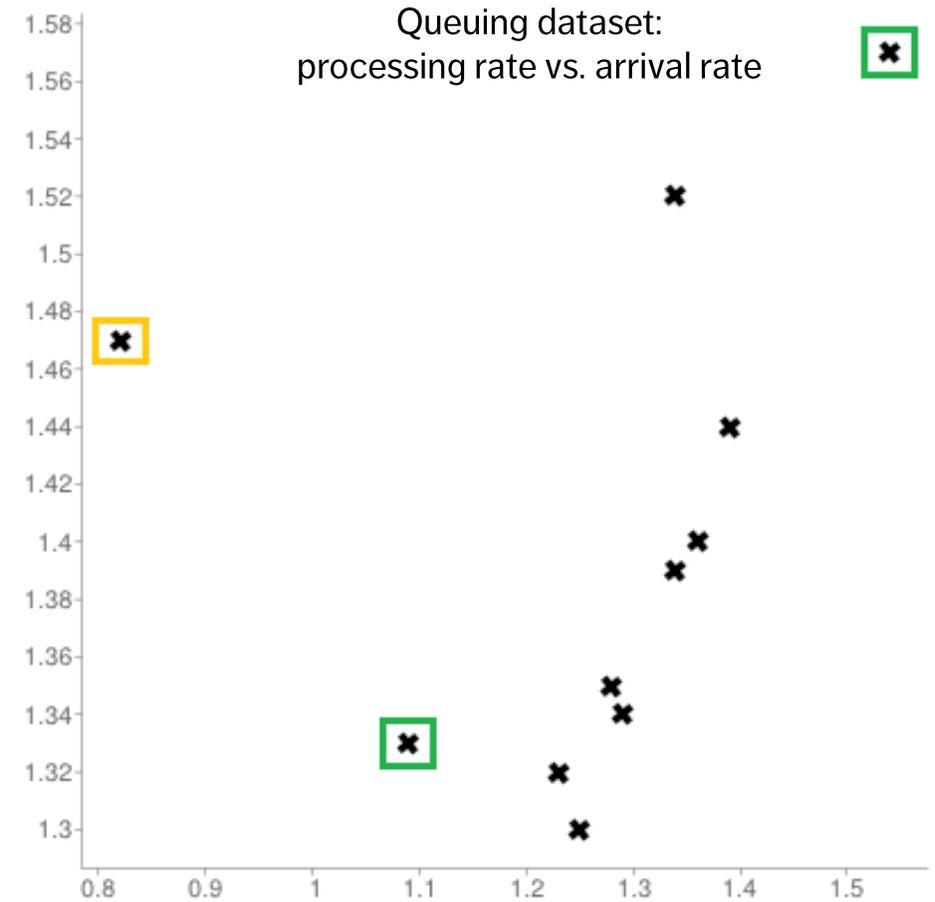
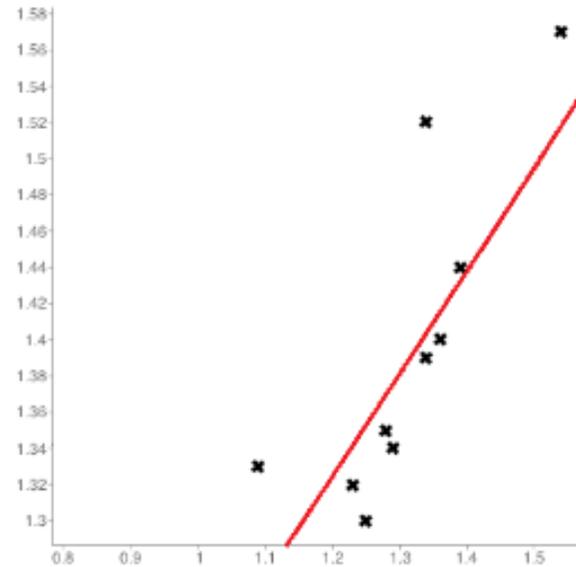
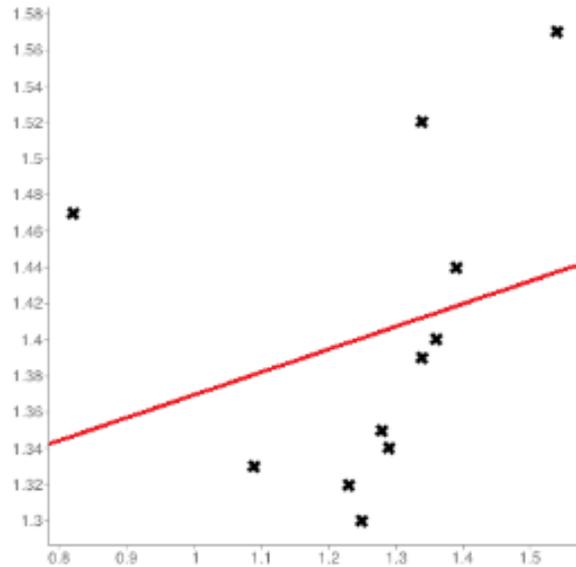


When influential observations are identified, **remedial measures** (such as data transformations) may be required to minimize their undue effects.



Outliers may be influential data points; influential data points need not be outliers (and *vice-versa*).

# INFLUENTIAL OBSERVATIONS



# ANOMALY DETECTION REMARKS

Identifying influential points is an **iterative process** as the various analyses must be run numerous times.

Fully automated identification and removal of anomalous observations is **NOT recommended**.

Use data transformations if the data is **NOT normally distributed**.

Whether an observation is an outlier or not depends on **various factors**; what observations end up being influential data points depends on the **specific analysis to be performed**.

---

# DATA ANALYSIS

OUTLIERS AND ANOMALOUS OBSERVATIONS