

# Data quality and validation

Part of the material Adapted from  
presentation prepared by Jon Wang ,  
Monika Kuffer  
ITC, University of Twente  
And Gis Gate

October 2022

The logo for IDEA MAP SUDAN is displayed in a large, semi-transparent white circle. The word 'IDEA' is in a large, bold, black font, with a vertical bar to its left containing three colored squares (yellow, blue, blue). The word 'MAP' is in a similar bold, black font, with a blue square to its left. The word 'SUDAN' is in a bold, black font below 'MAP'.

**IDEA**  
**MAP**  
**SUDAN**



African Population and  
Health Research Center



ULB



**nuffic**  
meet the world

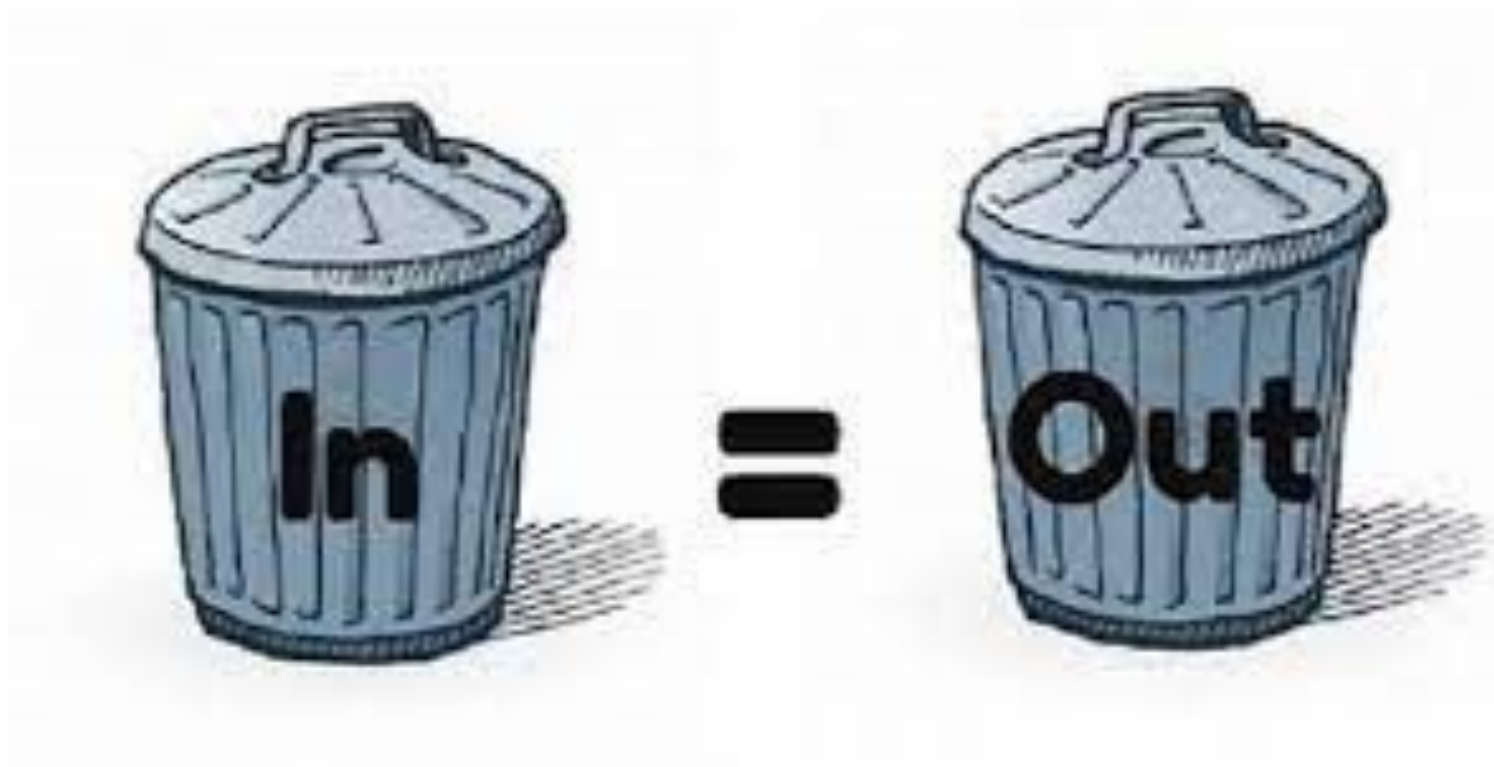
- Importance of data quality
- Definitions
- Topology
- Validation



## Data quality

Why data quality is important issue?

---



## Why data quality is important issue?

---

**Valid Data**



**Valid Analysis**



**Valid Result**

Moving from Accuracy precision uncertainty and Errors To Quality.

### □ **Accuracy :**

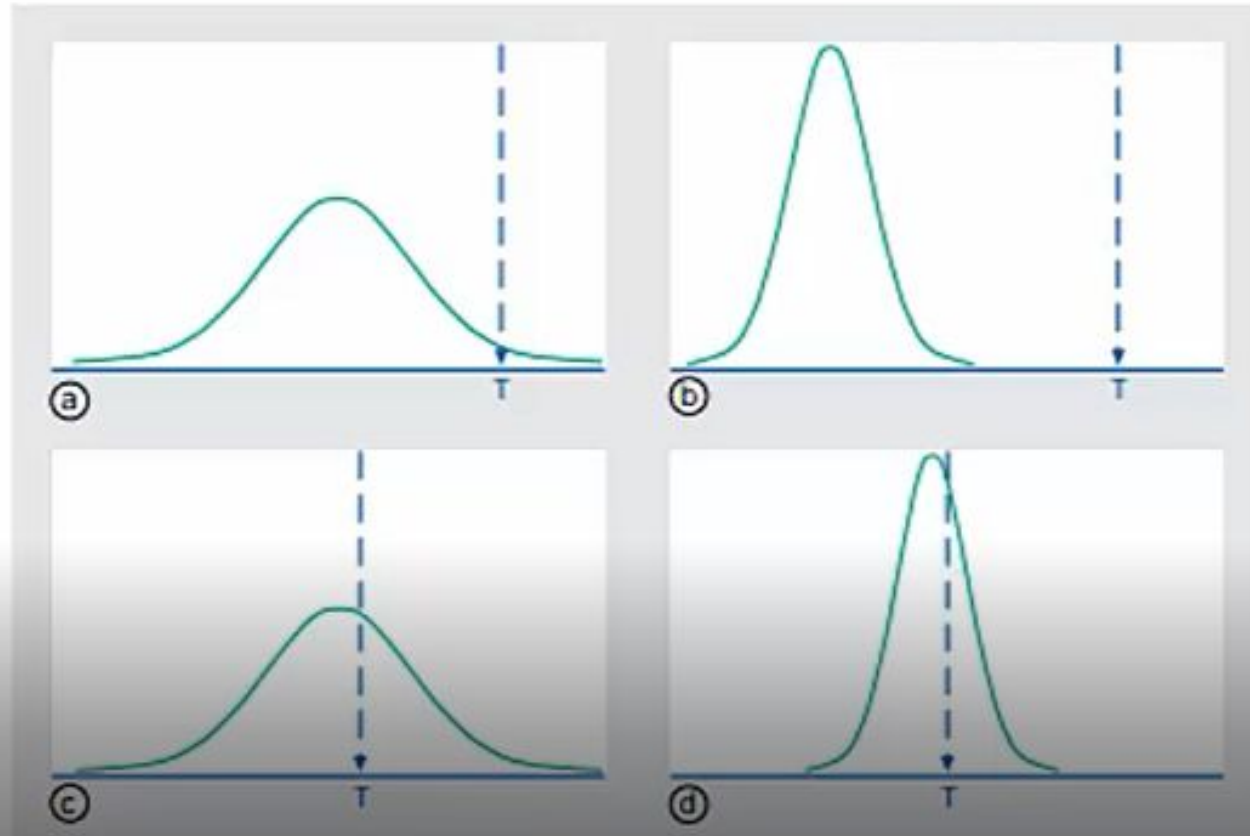
- Closeness of observation to it's true value
- Accurate measurement has a mean close to the true value

### □ Precision

- The smallest unit of measurements to which data can be recorded
  
- A precise measurements has small variance

# Data quality

## Basic Definition and Terminology



**Figure 8.50**

A measurement probability function and the underlying true value T: (a) bad accuracy and precision, (b) bad accuracy/good precision, (c) good accuracy/bad precision, and (d) good accuracy and precision.

### □ Error :

The different between measurement value and a true or theoretically correct value

### □ Type of errors:

- Human errors ( collecting, digitizing ,interpretation....)
- Instruments error ( systematic error not calibrated instrument)
- Random errors ( natural variation , small error)



### □ **Uncertainty:**

- Simply you are not sure about your measurements
- Uncertainty is not an error
- Errors usually originate from uncertainty

### □ **Data Quality :**

Characteristic of a product or service that bear on it's ability to satisfy stated and implied needs

- Quality is relative
- Quality doesn't mean excellence
- Quality is measured as fitness for purpose

# Data quality

Elements of spatial data quality :

---

- **Accuracy** ( positional, Thematic)
- **Lineage** (family tree of the data set)
- **Logical consistency** ( value should not be there )
- **Completeness** (you have the geometry but not have the names)

These and other data what expected to see in meta data

•

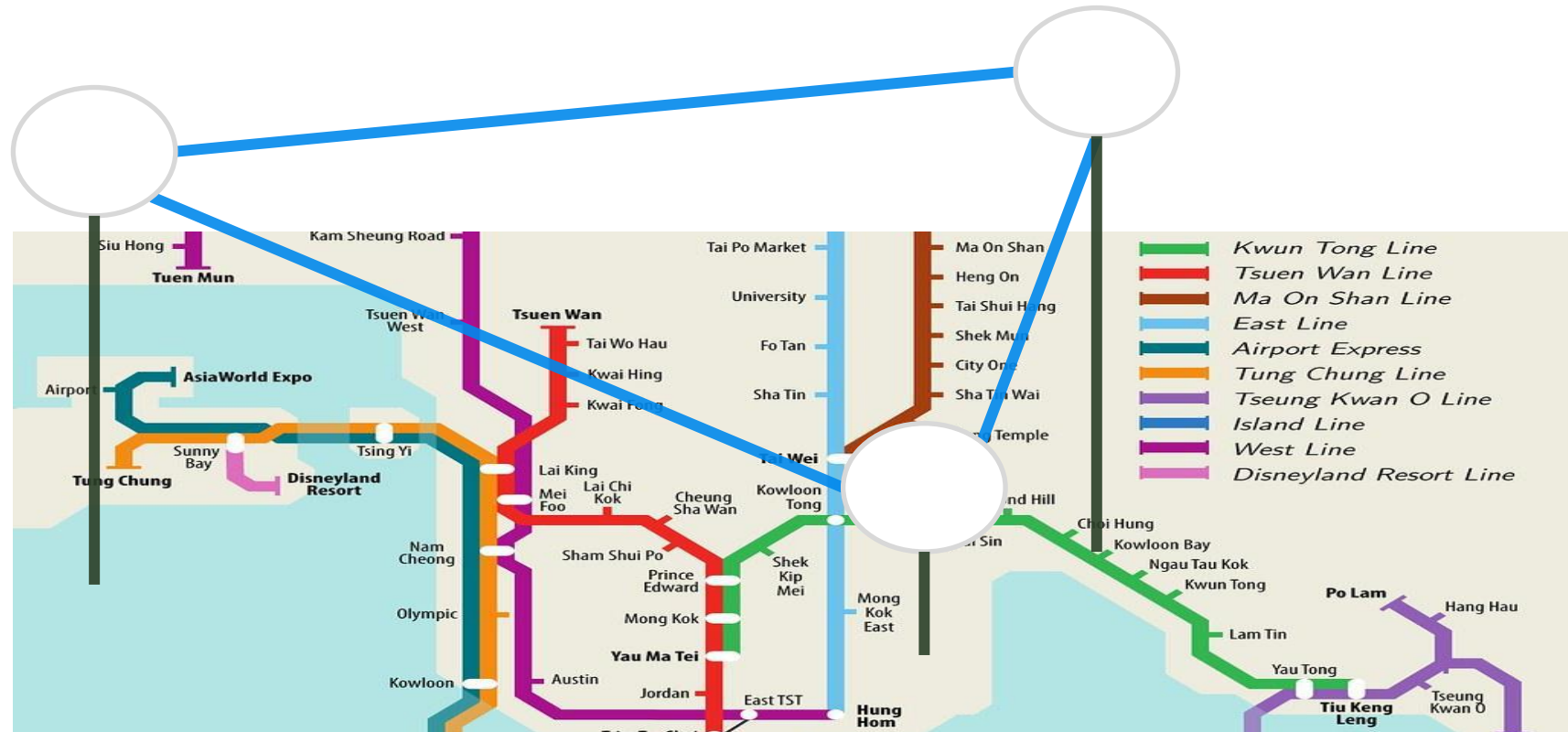
## How to deal with incompleteness or erroneous data in general

---

- Completeness is not always there so managing data gaps is essential in research
- There is no single approach to solve this problem but basic principles to keep in mind:
  - Be consistent especially if the goal is to compare
  - Depending on the nature of the data (estimating , deleting , ignoring) missing values might be an option
  - Always document what you did and be consistent

# What is a Topology?

- Essentially, topology refers to the relationship between things.
- In GIS, topology refers to the relationship between spatial objects.





□ What is topology?

In GIS relation between spatial objects

□ Why topology is important ?

Topology can guarantee that the acquired data is valid

—————> Valid analysis      —————> Valid results

### **Topology relationship:**

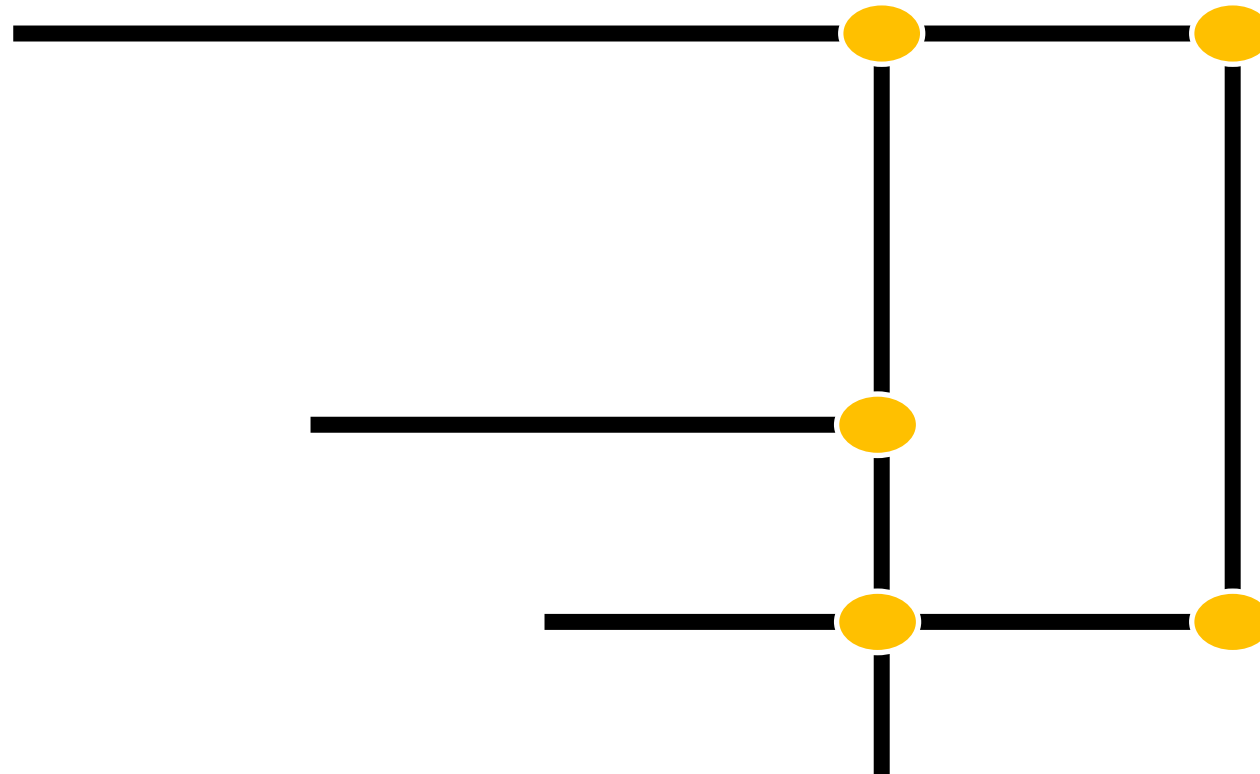
- Connectivity
- Adjacency
- Inclusion

# Topological Relationship

---

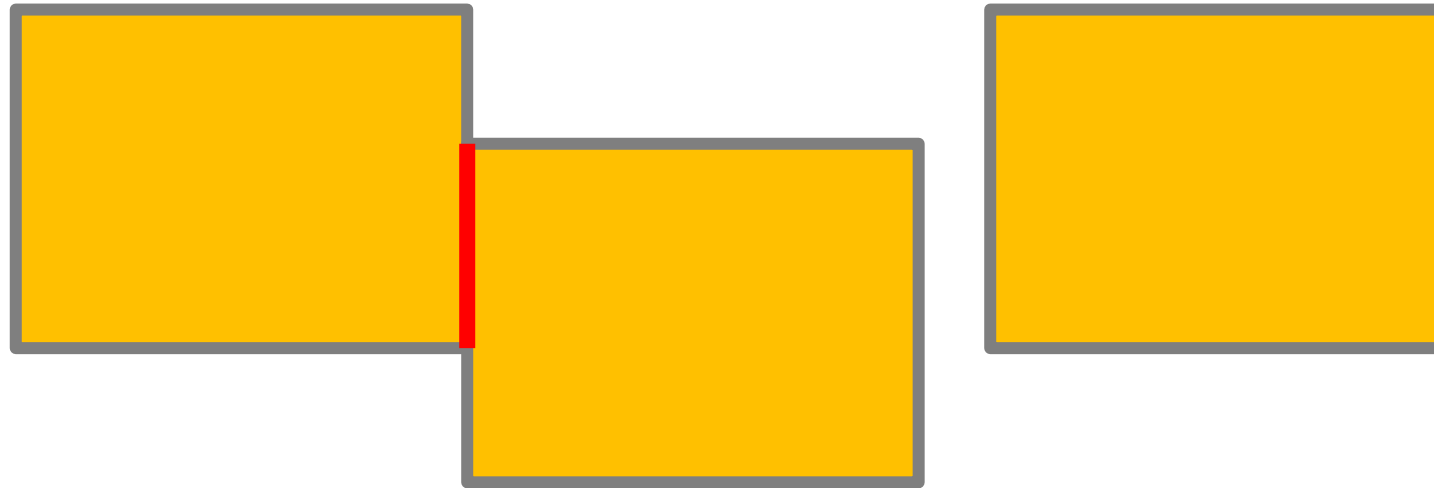
## □ Connectivity

Describes how lines are connected to each other to form a network



## ▣ Adjacency

Describes whether two areas are next to each other.

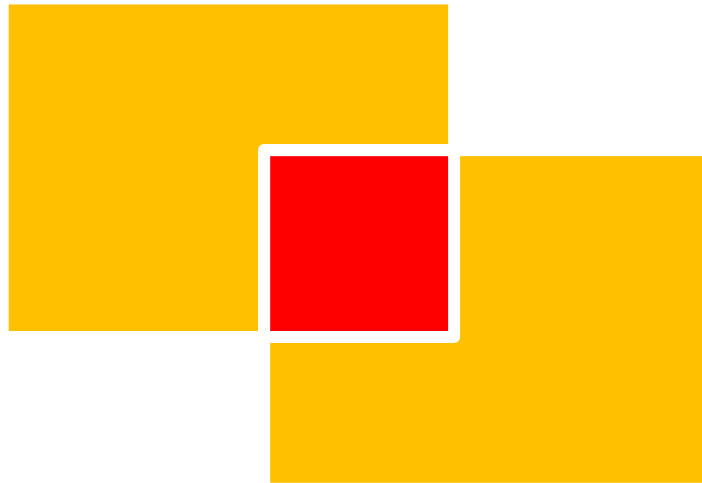


# Topological Relationships

---

## □ Inclusion

Describes whether two areas are nested.





## Topology Data Structure

---

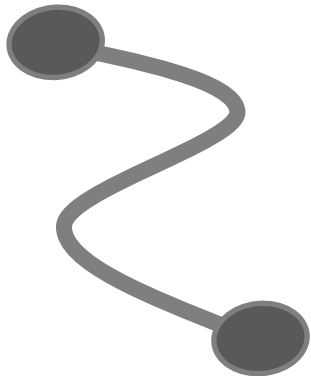
The most common topological data structure is the **Arc/Node** data model



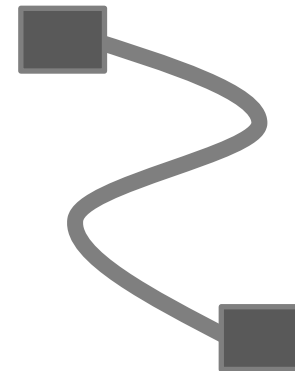
**Point**



**Node**



**Line**

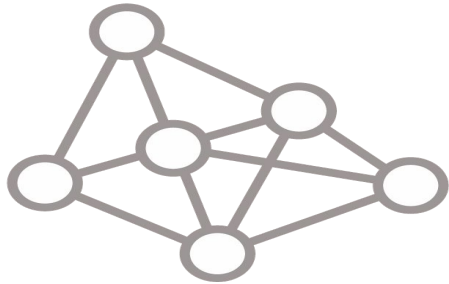


**Arc/Edge**

## How does topology work?

---

- Defining topology rules
- Validating data
- Resolving violation

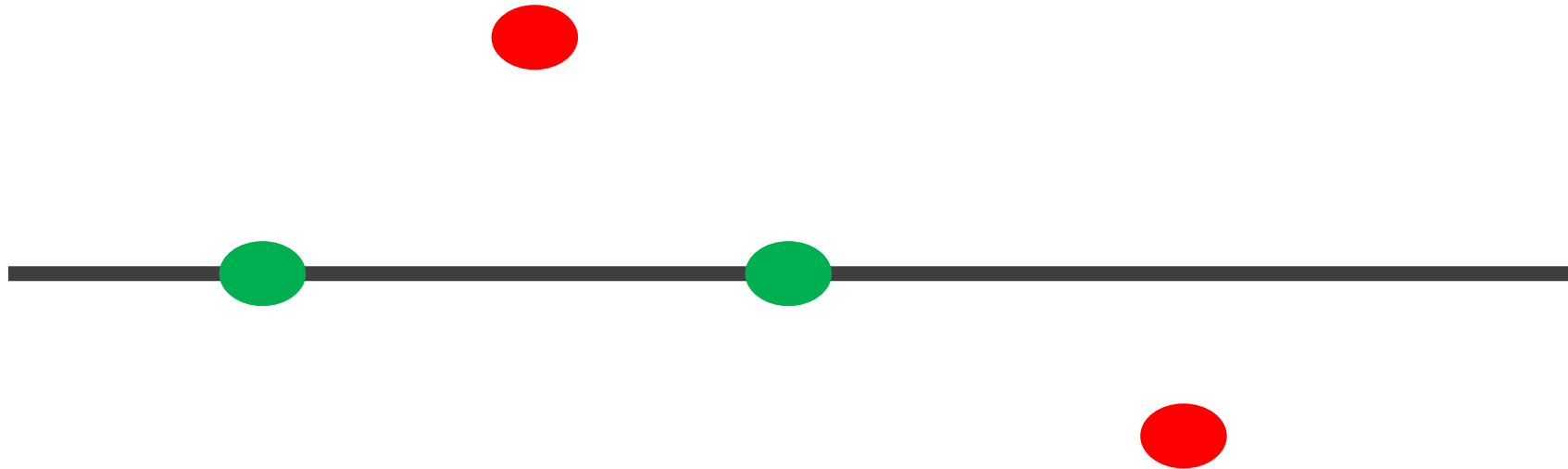


### **Topology Rules:**

- Define rules and constrains for the feature classes
- +30 different type of topology rules
- Applied in point, line and polygon feature
- Violation expressed by topology errors

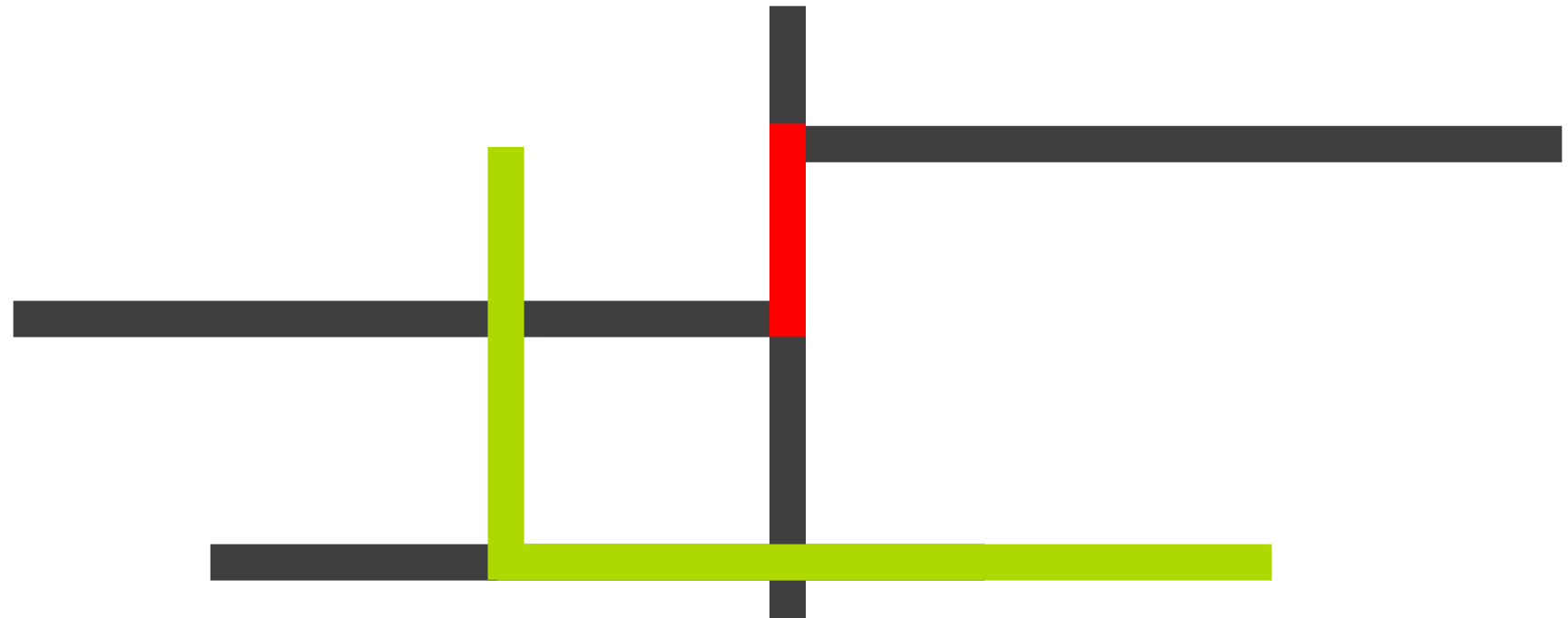
### □ Point Geometry

*Points Must Be Covered By Line*



## □ Line Geometry

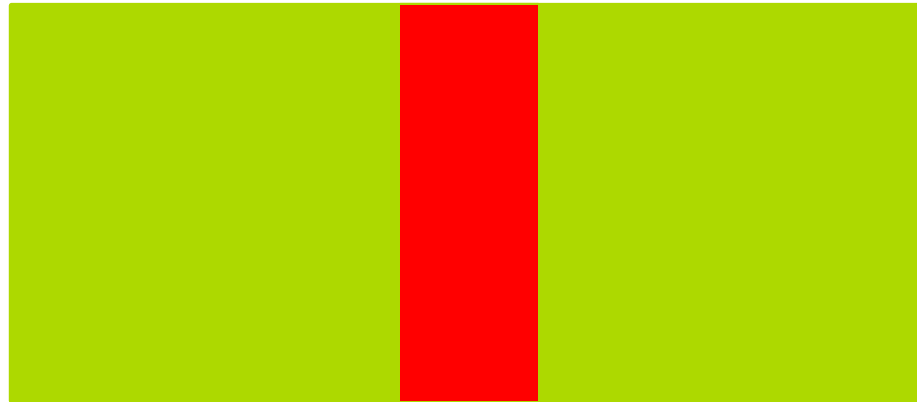
*Lines Must Not Overlap*





## □ Polygon Geometry

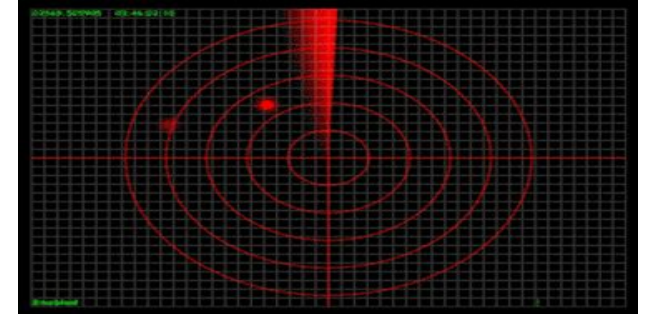
*Polygons Must Not Overlap*



# Validation

---

- It is the process of scanning data for rules violation Scanning the whole data or the visible data



## Topology layers:

- Point errors
- Line errors
- Area errors

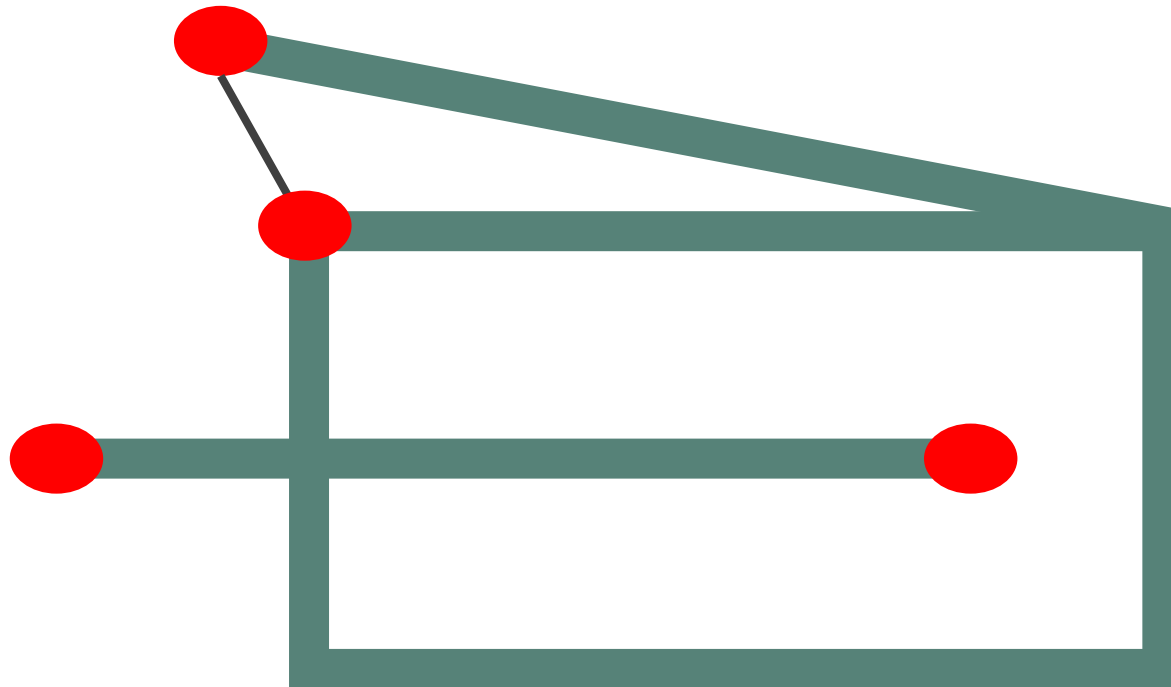


# Topology Solutions for Line Errors

---

## □ **Case study:** Lines must not have dangles

- Trim
- Extend
- Snap

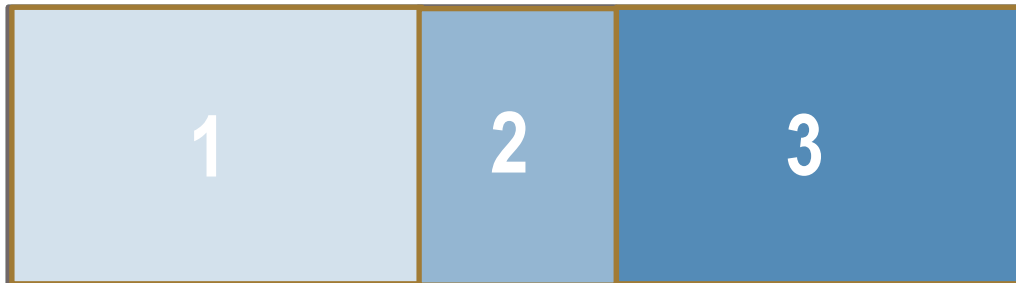


## Topology Solutions for Polygon Errors

---

### □ Case study: **Polygons must not overlap**

- Subtract
- Merge
- Create feature



- Topology is built on data base level
- Target layers should be determined along with their ranks
- Define Topology rules
- Validate data
- Fix errors

**The most important thing when dealing with errors is the consistency**

# Questions?

Selftest

---

