# DATA PREPROCESSING

Dr. Edi Surya Negara, M.Kom. Postgraduate Program, Informatics Engineering (S2) May 11, 2018





# **Chapter 1: Objectives**

After completing this chapter, students will be able to:

- Explain how to data cleaning.
- Explain how to data integration.
- Explain how to data transformation.

Reference: Tan et al: Introduction to data mining. Some slides are adopted from Tan et al.



## **Objectives and Benefits Data Preprocessing**

#### Preprocessing Steps in Data Mining



wis transformation -2, 32, 100, 99, 48 -+ -0.02, 0.32, 1.00, 0.59, 0.48



- Reduction of computing time.
- Can make the data value becomes smaller without changing the information it contains.





Data quality is measured from several aspects:

- Accuracy. The degree of conformity of a measure to a standard or a true value.
- **Completeness**. The degree to which all required measures are known.
- **Consistency**. The degree to which a set of measures are equivalent in across systems.
- Timeliness. Whether the data is uptudate.
- **Trusted**. The degree to which the measures conform to defined business rules or constraints.
- **Interpreted**. Whether the data can be understood.





# What is data cleaning?

Data cleaning is the process of detecting and correcting (or removing) corrupt or inaccurate records from a record set, table, or database and refers to identifying incomplete, incorrect, inaccurate or irrelevant parts of the data and then replacing, modifying, or deleting the dirty or coarse data.

Data cleaning includes:

- Fulfillment missing values.
- Smooth data noisy.
- Identify outliers.
- Eliminating inconsistencies.



## Attribute value error could be due

- Error instrument data collection
- Data entry error
- Data delivery error
- Technological limitations
- Inconsistent naming conventions
- Records are repeated
- Incomplete data
- Inconsistent data





## **Data Transformation**

Data Transformation Goals:



- Overcome challenge of variable underlying data set structures through:
  - Creating a uniform, integrated data set that allows for timely and easily accessible reports
  - Integrating data needs according to a central schema



### Several ways to data transformation

- Centering
- Normalization z-core
- Scaling



(1

# Centering

- Based on fitting a distribution to the data
- Distance function between distributions

Centering formulation:

$$\widehat{X} = X - \overline{X}$$

 $\widehat{X}$  = vector results after centering X = vector of the original column  $\overline{X}$  = average of the columns in question.





**Scatter Matrix** is the distance between attributes  $X_1$  with  $X_2$ , or  $X_1$  with  $X_3$  and so on.

$$Scatter = \widehat{X}^T \widehat{X}$$
(2)

**Kovarian Matrix** is the distance between attributes  $X_1$  with  $X_2$ , or  $X_1$  with  $X_3$  and so o, but divided by each entry in the matrix by the amount of data m - 1.

$$xovarian = \frac{\widehat{X}^T \widehat{X}}{m-1}$$
(3)



#### Normalization z-core

**Normalization** dividing each data that has been cast with the standard deviation of the attribute in question.

$$\widehat{X} = \frac{X - \overline{X}}{\sigma_{X}} \tag{4}$$

Data after this standardization will have a mean of = 0 and variance  $= 1 \label{eq:constraint}$ 



### Scaling

Scaling procedure to change the data so that it is in a certain scale.

$$\widehat{X} = \frac{X - X_{min}}{X_{max} - X_{min}} * (BA - BB) + BB$$
(5)





#### Data Transformation: Example

```
FATFACHING/Data Mining dan Rig Data Analytics/Data-Mining-and-Rig-Data-Analytics-Rock-master/Data-Mining-and-Rig-Data-Analytics-Rock-master/3.4 Transformasi Data.nv - Notenad++
File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window
        client] ov 🛛 🚍 server1 ov 🖓 🚍 chat server ov 🖓 🚍 plot stock market ov 🖓 🚍 34 Transformasi Data.ov 🕅
    #import numpy
    import numpy as np
    #mendefinisikan matrik x, bebas, untuk contoh saja
    x = np.matrix([[-20,23,5],[4,-8,15]])
  7 #centering (data asli dikurangi nilai reratanya)
  8 #sesuai Persamaan 1
  9 x centering = x - x.mean()
 10 print("x centering: ", x centering, "\n")
 12 #standarisasi sesuai Persamaan 4
 13 x standarisasi = (x-x.mean()) / x.std()
    print("x standarisasi: ", x standarisasi)
     print("mean x standarisasi: ", x standarisasi.mean())
 16 print("varian x standarisasi: ", x standarisasi.var(), "\n")
 18 #scaling ke range 0-1, sesuai Persamaan 5
 19 BA = 1; BB = 0 #BA=batas atas, BB = batas bawah
 20 x scaling = (x - x.min()) / (x.max()-x.min()) * (BA-BB) + BB
 21 print("x scaling: ", x scaling)
```





#### Data Transformation: Output

