

Chapter 3: Data Preprocessing

Dong-Kyu Chae

**PI of the Data Intelligence Lab @HYU
Department of Computer Science & Data Science
Hanyang University**



Contents: Major Tasks in Data Preprocessing

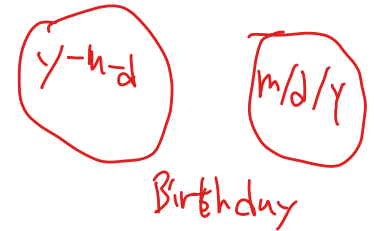
- ☐ **Data Cleaning**
- ☐ **Data Integration**
- ☐ **Data Reduction**
- ☐ **Data Transformation and Data Discretization**



Why Pre-process the Data?

❑ A multidimensional view

- ❑ **Accuracy:** your data mining/machine learning results are not accurate, despite multiple trials...
- ❑ **Completeness:** not recorded, unavailable, ...
- ❑ **Consistency:** some modified but some not, ...
- ❑ **Timeliness:** timely updated? *data / t*
- ❑ **Believability:** how trustable the data are?
- ❑ **Interpretability:** how easily the data can be understood?



preprocessing이 필요한 이유



Overview

❑ Data cleaning

dirty하게 만들기

- ❑ Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies

❑ Data integration

Scatter되어 있는걸 integrate

- ❑ Integration of multiple databases, data cubes, or files

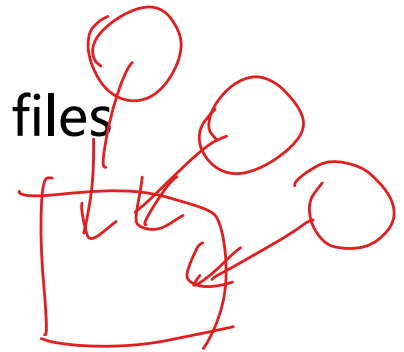
❑ Data reduction

- ❑ Dimensionality reduction
- ❑ Numerosity reduction
- ❑ Data compression

여러개로 data를
묶음

❑ Data transformation and data discretization

- ❑ Normalization
- ❑ Concept hierarchy generation





Data Cleaning

- ❑ **Data in the Real World Is Dirty:** Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, transmission error 수집된 데이터가 videost-terrible
- ❑ incomplete: lacking feature values, lacking certain features of interest, or containing only aggregate data
 - e.g., *Occupation*= " " (missing data) missing value
- ❑ noisy: containing noise, errors, or outliers
 - e.g., *Salary*= "-10" (an error)
- ❑ inconsistent: containing discrepancies in codes or names, e.g.,
 - *Age*= "42" , *Birthday*= "03/07/2000" inconsistent
 - In some DBs, *rating* is "1, 2, 3" , but some other DBS, "A, B, C"
 - discrepancy between duplicate records

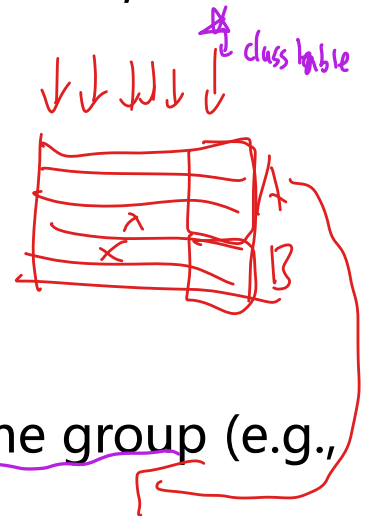


Missing Data

missing value를 ignore
한다.

- ❑ **Remove the object:** usually done when class label is missing—
not effective when the % of missing values is large
- ❑ **Fill in the missing value manually:** might be accurate, but
tedious + infeasible
cost ↑
- ❑ **Fill in it automatically with**
 - ❑ a simple constant (default value, or “unknown”)
 - ❑ the feature mean
 - ❑ the feature mean for all samples belonging to the same group (e.g.,
same class, same cluster, etc...)
 - ❑ **the inferred value:** such as based on some regression or
classification model

class label이 X이면
전체를 삭제하거나 다른
feature가 X이면 삭제X



missing value를
ML model로 predict
feature로
Dollars



Noisy & Inconsistent Data

❑ **Noise** : random error or variance in a measured feature

❑ Mainly due to faulty data collection instruments

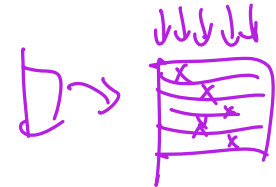
❑ Noisy data is often expressed as an **outlier**

▪ Outlier detection -> delete outliers -> find missing values

❑ Thus we can apply an **outlier detection** method (will learn it later)

salary 90000
= -10 to 50 110000
outlier

noise data 2 remove



❑ **Inconsistent data**

❑ Age= "42" , but Birthday= "03/07/2000"

❑ For a duplicate records, one name is "cm " but the other is "inch"
group을 정하고 → convert 단위. 'inch' → cm

❑ Human inspection will be needed ~> 너무 많이 필요할 것 같다.

▪ Computer performs outlier detection, then human will inspect it



Contents

❑ Data Preprocessing: An Overview

- ❑ Data Quality
- ❑ Major Tasks in Data Preprocessing

❑ Data Cleaning

❑ Data Integration

❑ Data Reduction

❑ Data Transformation and Data Discretization

❑ Summary



Data Integration

❑ Data integration:

- ❑ Combines multiple datasets from multiple sources into a coherent store

❑ Schema integration: e.g., $A.cust-id \equiv B.cust-\#$

- ❑ Integrate metadata from different sources

이름이 다른 것
meaning이 같은 것



❑ Detecting and resolving data value conflicts

- ❑ For the same real world entity, feature values from different sources are different
- ❑ e.g., cm vs. inch, meter vs. mile

401에서 mile
301에서 km
인정 conflict

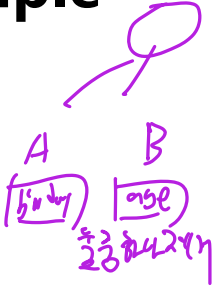


Handling Redundancy in Data Integration

❑ Redundant data occur often when integration of multiple databases

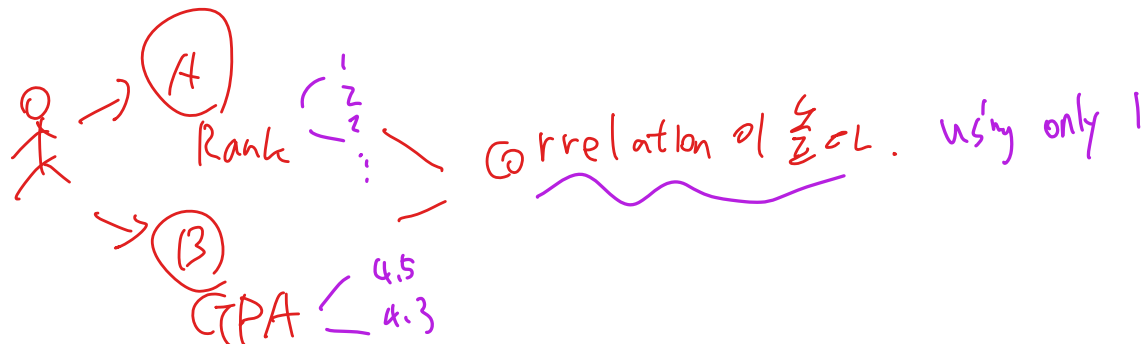
- ❑ Derivable data : One feature may be a “derived” feature in another table, e.g., birthdate vs. age

redundant 중복하나제거



❑ Redundant features can be automatically detected by correlation analysis and covariance analysis

❑ Reducing/avoiding redundancies and inconsistencies improves mining speed and quality



Correlation Analysis (Nominal Features)

❑ We want to know that “like_science_fiction” and “play_chess” are **correlated**

	Play chess (Y)	Not play chess (N)	Sum (row)
Like science fiction (Y)	250(90)	200(360)	450 1:4
Not like science fiction (N)	50(210)	1000(840)	1050 1:4
Sum(col.)	300	1200	1500 1:4

❑ **X² (chi-square) test**

$$\chi^2 = \sum_{\text{each cell}} \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

- ❑ The larger the X² value, the more likely the features are corelated
- ❑ The cells that contribute the most to the X² value are those whose actual count is very different from the expected count
- ❑ Expected value is estimated under the independence assumption

Handwritten notes on the slide:

- Red arrows pointing to the expected values in parentheses: "expected values"
- Red text: "independent 가정" (independence assumption)
- Red text: "노오려서 210" (likely 210)
- Red arrows pointing to the 1:4 ratios: "1:4"
- Red text: "↓ ~ independent", "↑ ~ correlated"

Correlation Analysis (Nominal Features)

	Play chess	Not play chess	Sum (row)
Like science fiction	250(90)	200(360)	450
Not like science fiction	50(210)	1000(840)	1050
Sum(col.)	300	1200	1500

❑ χ^2 (chi-square) calculation

❑ Numbers in parenthesis are expected counts calculated based on the data distribution in the two categories

$$\chi^2 = \frac{(250 - 90)^2}{90} + \frac{(50 - 210)^2}{210} + \frac{(200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840} = 507.93$$

chi-square table
여기에
값을
넣는다
correlation
표는 21

❑ It shows that like_science_fiction and play_chess are correlated in the group



Correlation Analysis (Numeric Features)

- Correlation coefficient (also called **Pearson's correlation coefficient, PCC**) among features A and B:

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \overline{A})(b_i - \overline{B})}{(n-1)\sigma_A\sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n\overline{A}\overline{B}}{(n-1)\sigma_A\sigma_B}$$

Handwritten notes:
- Above \overline{A} and \overline{B} : average
- Below $(n-1)$: number of data points

n is the number of data, \overline{B} and \overline{A} are the respective means of A and B, σ_A and σ_B are the respective standard deviation of A and B, and $\sum(a_i b_i)$ is the sum of the AB cross-product.

- If $r_{A,B} > 0$: A and B are **positively correlated**

- A's values increase as B's). The higher, the stronger correlation

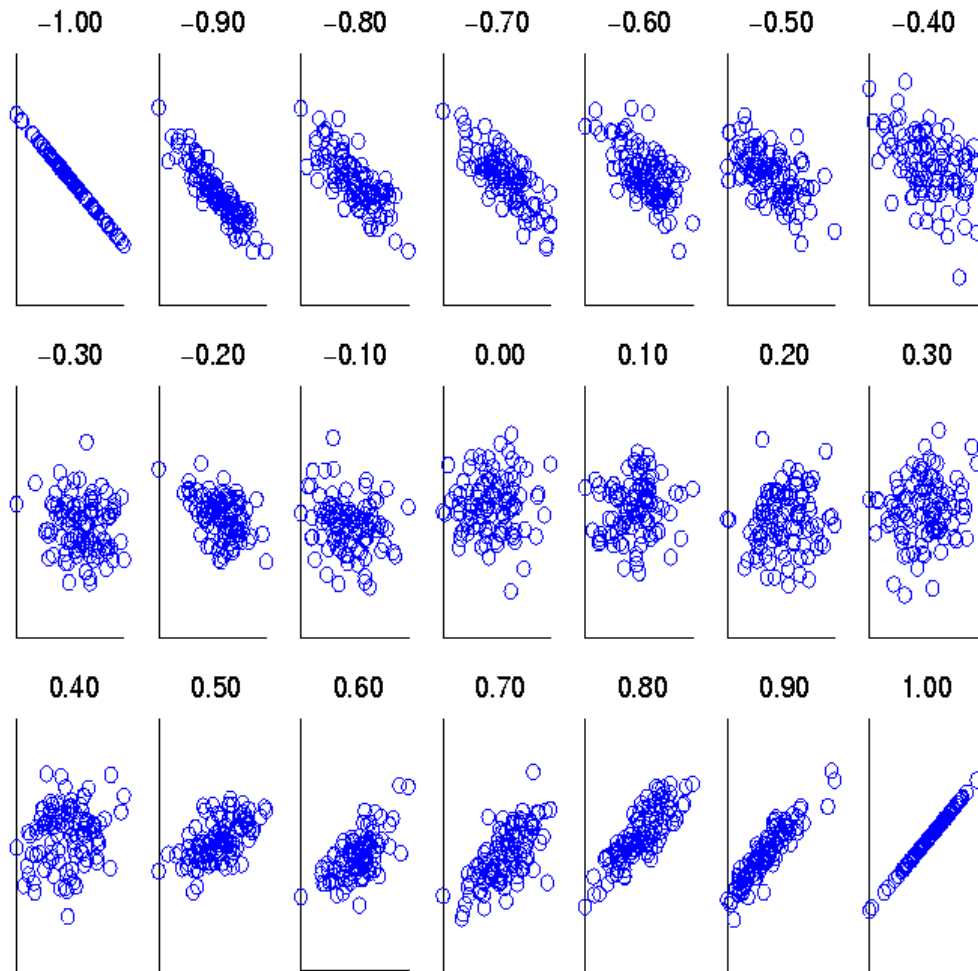
- $r_{A,B} = 0$: **independent** *completely independent*

- $r_{AB} < 0$: **negatively correlated**

Handwritten note: A가 증가하면 B도 tend to be increasing



Visually Evaluating Correlation



Scatter plots showing the correlation from -1 to 1 .

Correlation(상관관계) does not imply causality(인과관계)
=> “# of hospitals” and “# of car-theft” in a city are correlated. However, both may be causally linked to another feature:
population



Covariance (Numeric Features)

- **Covariance is similar to correlation**

$$Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

Correlation coefficient: $r_{A,B} = \frac{Cov(A, B)}{\sigma_A \sigma_B}$

where n is the number of data, \bar{A} and \bar{B} are the respective mean or expected values of A and B , σ_A and σ_B are the respective standard deviation of A and B .

- **Positive covariance:** If $Cov_{A,B} > 0$, then A and B both tend to be larger than their expected values.

- **Negative covariance:** If $Cov_{A,B} < 0$ then if A is larger than its expected value, B is likely to be smaller than its expected value.

- **Independence:** $Cov_{A,B} = 0$

specific range
that
↓
normalize
cov



Covariance: An Example

skip

$$Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

- It can be simplified in computation as

$$Cov(A, B) = E(A \cdot B) - \bar{A}\bar{B}$$

- Suppose two stocks A and B have the following values in one week: (2, 5), (3, 8), (5, 10), (4, 11), (6, 14).

- Question: If the stocks are affected by the same industry trends, will their prices rise or fall together?

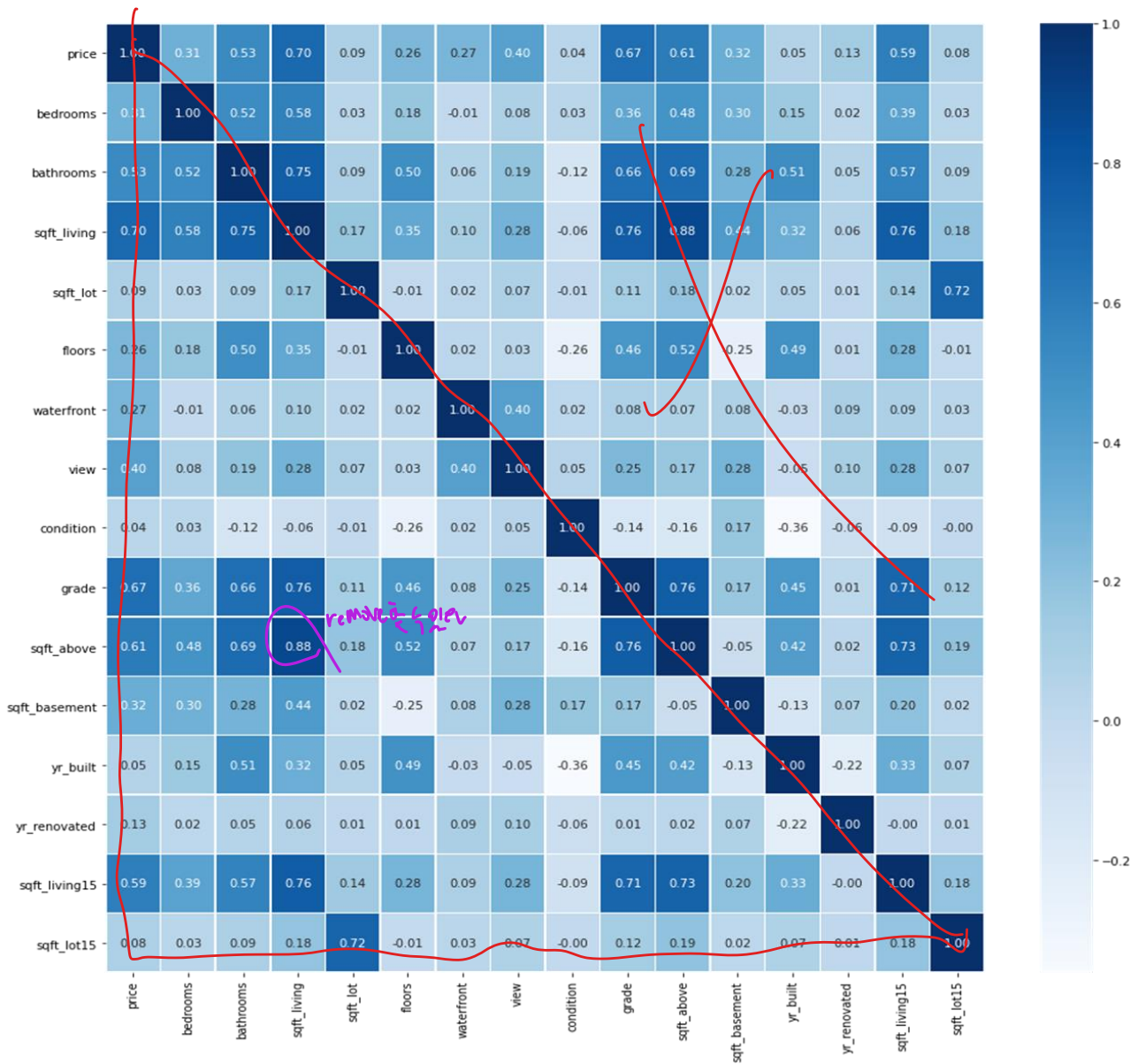
- $E(A) = (2 + 3 + 5 + 4 + 6) / 5 = 20 / 5 = 4$

- $E(B) = (5 + 8 + 10 + 11 + 14) / 5 = 48 / 5 = 9.6$

- $Cov(A, B) = (2 \times 5 + 3 \times 8 + 5 \times 10 + 4 \times 11 + 6 \times 14) / 5 - 4 \times 9.6 = 4$

- Thus, A and B rise together since $Cov(A, B) > 0$.


Covariance/Correlation Matrix Visualization



highly
correlation
이런
특징을
선택



Contents

- ❑ **Data Preprocessing: An Overview**
 - ❑ Data Quality
 - ❑ Major Tasks in Data Preprocessing
- ❑ **Data Cleaning**
- ❑ **Data Integration**
- ❑ **Data Reduction** 
- ❑ **Data Transformation and Data Discretization**
- ❑ **Summary**

Thank You