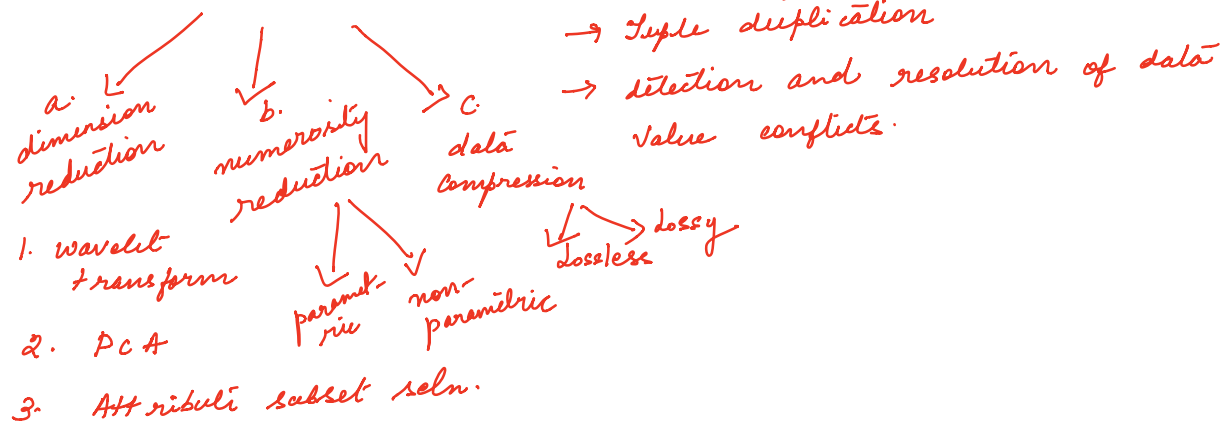


19/01/21 Dec-09 → Data Integration

### Data Reduction



### ① Entity Identification Problem:

Eg: marks: - 1. age  
cgpa (1 - 10)  
Grade (A, A+, ...)

Focus on the semantics of data i.e. maintain meta data (for each attribute)

This includes name, meaning, data type, range of values permitted for the attribute and null rules for handling blank, zero or null values.

→ The meta data can be used to avoid errors in schema integration.

### 2) Redundancy and Correlation analysis -

It is an imp. issue in data integration.

→ The redundancies can be detected using correlation analysis.

→ This measures, given two attributes, how strongly one attribute implies the other, based on available data.

for nominal attribute:  $\chi^2$  (chi-square test)

for numerical: → Pearson Correlation Coefficient  
→ Covariance Test.

a) Chi-square Test:

Diagram showing a contingency table structure:

	$A_1$	$A_2$	...	$A_c$
$B_1$	$(A_i B_j)$			
$B_2$				
$\vdots$				
$B_r$				

Overall dimensions:  $r \times c$

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E_{ij} = \frac{\text{Count}(A=a_i) \times \text{Count}(B=b_j)}{N}$$

Null

Hyp: A & B are independent

$$\chi^2 \rightarrow 25 \checkmark$$

$$\alpha \rightarrow 0.1 \checkmark$$

$$0.01$$

$$0.001 \checkmark$$

$$0.0001$$

$$\text{DOF} \rightarrow (r-1) \times (c-1)$$

$$(2-1) \times (2-1)$$

$$\frac{10}{0.001} = 10000$$

1 → 10 ✓

$O_{ij}$  → observed frequency (actual count) of joint event  $(A_i B_j)$

$E_{ij}$  → expected " of  $(A_i B_j)$

$N$  = no. of data tuples

$\text{Count}(A=a_i)$  is the no. of tuples having value  $a_i$  for A.

$\text{Count}(B=b_j)$  " " " " " " "  $b_j$  for B.

	Male	female	
fiction	250 (90)	200	→ 450
non-fiction	50	1000	→ 1050
Total	300	1200	(1500)

$\chi^2$  test hypothesis (Null hyp.) assumes A & B are independent (gender & Types of book are independent)

$$e_{11} = \frac{\text{Count (male)} \times \text{Count (fiction)}}{N}$$

$$= \frac{300 \times 450}{7500} = 90$$

$$e_{12} = \frac{1200 \times 450}{1500} = 360$$

$$e_{21} = \frac{300 \times 1050}{1500} = 210$$

$$e_{22} = \frac{1200 \times 1050}{1500} = 840$$

$$\begin{aligned} \chi^2 &= \frac{(250 - 90)^2}{90} + \frac{(200 - 360)^2}{360} + \frac{(50 - 210)^2}{210} \\ &\quad + \frac{(1000 - 840)^2}{840} \\ &= \underline{\underline{507.93}} \end{aligned}$$

$$DOF = (r-1) \times (c-1) = (2-1) \times (2-1) = 01$$

For 1 DOF, the  $\chi^2$  value needed to reject the hypothesis at 0.001 significance level is 10.828 (refer  $\chi^2$  distribution Table). Since, the computed value is above this, we reject the hypothesis that gender and book reading are independent.

Therefore, we conclude that the two attributes are dependent or strongly correlated to each other.

b) Pearson Correlation Coefficient:

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

$$-1 \leq r_{A,B} \leq +1$$

If  $r_{A,B}$  is  $> 0$ , then A and B are positively correlated  
(ie  $A \uparrow \quad B \uparrow$ )

If  $r_{A,B} = 0$ , then A and B are independent and  
there is no correlation b/w them.

If  $r_{A,B} < 0$ , then A and B are negatively correlated  
i.e  $A \uparrow \quad B \downarrow$

Eg:

serial no.	Age (x)	Glucose level (y)	xy	$x^2$	$y^2$
1	43	99	4257	1849	9801
2	21	65	1365	441	4225
3	25	79	1975	625	6241
4	42	75	3150	1764	5625
5	57	87	4959	3249	7569
6	59	81	4779	3481	6561
$\Sigma$	247	486	20485	11409	40022

$r = \underline{\underline{0.5298}}$  (This means 52.98% variables have a moderate positive correlation)

c) Covariance of numeric data:  
(Next class)