

21/01/21 Dec-10

Covariance :

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

for sample covariance :  $\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$

Relation b/w correlation  
& covariance :

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

→ Covariance is a measure of relationship b/w 2 random var.  
→ This measure evaluates to what extent the variables change together.

\* Positive covariance: indicates that 2 variables tend to move in same direction

\* Negative covariance: reveals that 2 variable move in opposite direction.

where :  $x_i$  : values of x variable

$y_i$  : " " y variable

$\bar{x}$  : mean / avg. of x var

$\bar{y}$  : mean / avg " y var.

$\sigma_x$  : std. dev. of x

$\sigma_y$  : " " " y.

data	$x$ TCS	$y$ Infosys	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$a * b$
2013	1692	68	-352.8	-41.20	14,535.36
2014	1978	102	-66.8	-7.2	480.96
2015	1884	110	-160.8	0.8	-128.64
2016	2151	112	106.2	2.8	297.36
2017	2519	154	474.2	44.8	21,244.16

$$\text{Mean } (TCS) : \frac{1692 + \dots + 2519}{5} = 2044.8$$

$$\text{Mean } (Infosys) \bar{y} : \frac{68 + 102 + \dots + 154}{5} = 109.2$$

$$\text{Cov}(x, y) : \frac{36429.20}{5} = \underline{\underline{7285.4}}$$

It can be concluded that, price of stock of TCS & Infosys move in same direction.

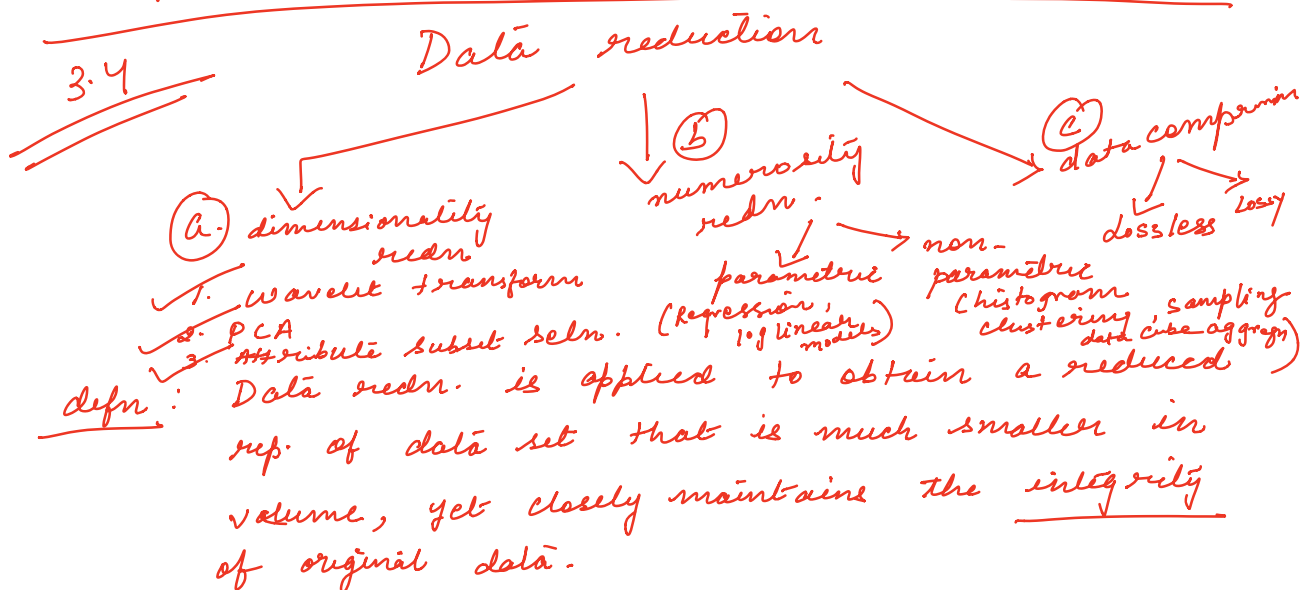
3  $\Rightarrow$  Tuple duplication:

Redundant tuples are identified and removed during integration

4 Data value conflict detection and resolution:

for Eg: wt. in metric unit  
wt in British Imperial unit

$\rightarrow$  The semantic integrity and str. of data are great challenge during data integration.



1) dimensionality reduction :

(Reducing attributes across dimension)

1) wavelet transform :

DWT is a linear signal processing technique that is when applied over a data vector  $X$ , transforms it to a numerically different vector  $X'$  of wavelet coefficients.

→ Fetching the strongest wavelet coefficients helps in lower dimension mapping.

$$X = \begin{pmatrix} 1 & \dots & 1000 \\ 0.5 & 0.1 \end{pmatrix} \xrightarrow{\text{DWT}} X' = \begin{pmatrix} 1 & \dots & 1000 \\ 0.225 & 0.618 \end{pmatrix}$$

$$X' = \begin{pmatrix} \dots & 100 \end{pmatrix}$$

$$\begin{pmatrix} 1 \\ \vdots \\ 100 \end{pmatrix} \xrightarrow{100 \times 1000} \downarrow \xrightarrow{100 \times 100} \text{DWT}$$

2) PCA :

$$A = \begin{pmatrix} \dots \end{pmatrix} \quad 100 \times 1000$$

$$\text{Cov. matrix} = A^T A : \begin{pmatrix} \dots \end{pmatrix}_{1000 \times 1000} \quad \begin{pmatrix} \dots \end{pmatrix}_{100 \times 1000}$$

$$: \begin{pmatrix} \end{pmatrix} \quad 1000 \times 1000$$

princomp (I) →

arranged in  
descending  
order

$$\begin{matrix} \begin{bmatrix} ev_1 \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_{1000} \end{bmatrix} \\ \swarrow \\ \begin{pmatrix} \lambda_1 & \lambda_2 & \lambda_3 & \dots & \lambda_{1000} \\ ev_1 & ev_2 & ev_3 & \dots & ev_{1000} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{pmatrix} \end{matrix} \quad 1000 \times 1000$$

$$\text{if } \left( \text{cum var: } \frac{\lambda_i}{\sum \lambda_i} > 0.99 \right)$$

$$\frac{\lambda_1 + \sqrt{\lambda_2 + \lambda_3}}{\sum \lambda_i} > 0.99 \quad ?$$

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_{100}}{\sum \lambda_i} > 0.99$$

$$I' = \begin{pmatrix} \lambda_1 & \dots & \lambda_{100} \\ ev_1 & & ev_{100} \\ \vdots & & \vdots \end{pmatrix} \quad 1000 \times 100$$

low. dim<sup>n</sup>. mapping:

$$\begin{pmatrix} \end{pmatrix}_{100 \times 1000} \times (I')_{1000 \times 100}$$

$$= \begin{pmatrix} \quad \end{pmatrix}_{100 \times 100}$$

3) Attribute seln:

$$f.v = \langle \check{f}_1 \text{ } \cancel{f_2} \text{ } \cancel{f_3} \text{ } \cancel{f_4} \text{ } \check{f}_5 \rangle$$

$$\{ \}$$

$$\{ f_1 \} \xrightarrow{ML} 90\%$$

$$\{ f_1, f_2 \} \xrightarrow{ML} 85\%$$

$$\{ f_1, f_3 \} \rightarrow 95\%$$

$$\{ f_1, f_3, f_4 \} \rightarrow 95\%$$

$$\{ f_1, f_3, f_5 \} \rightarrow 98\% \quad \checkmark$$

ii) backward seln:

$$\{ \cancel{f_1} \text{ } f_2 \text{ } f_3 \text{ } f_4 \text{ } f_5 \} \rightarrow 95\%$$

$$\{ \cancel{f_2} \text{ } f_3 \text{ } f_4 \text{ } f_5 \} \rightarrow 95\% \quad \cancel{90\%}$$

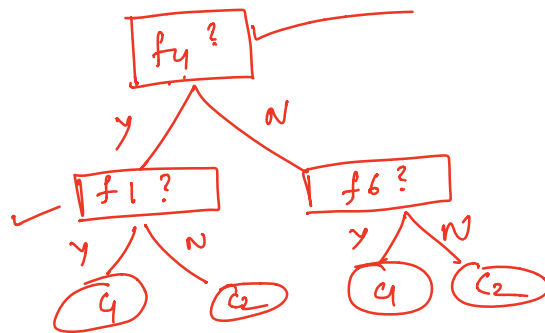
$$f_2 \text{ } \cancel{f_3} \text{ } f_4 \text{ } f_5 \rightarrow 95\% \quad 96\%$$

$$\{ f_1 \text{ } f_3 \text{ } f_5 \} \rightarrow 98\%$$

(iii) Combination of forward seln. & backward elimination

(iv) Decision Tree Induction :

Initial:  $\{ f_1 \ f_2 \ f_3 \ f_4 \ f_5 \ f_6 \}$



$\{ f_4 \ f_1 \ f_6 \}$