

Techniques : Data Reduction

- ① Data cube aggregation
- ② Attribute subset selection
- ③ Dimensionality Reduction
- ④ Numerosity reduction
- ⑤ Discretization and concept hierarchy generation

① Data cube aggregation

2017	Product	sales		Year	sales	
	H1	500	\Rightarrow	2017	700	} storing data year wise apply aggregation
	H2	200		2018	600	

2018	Product	sales
	H1	300
	H2	300

a_1	a_2	a_3	a_4	a_5	\Rightarrow	a_2	a_4	a_5
-------	-------	-------	-------	-------	---------------	-------	-------	-------

② Attribute subset selection

irrelevant att. may be discarded/removed/original.

Redundant/ a_k Sr. No. Roll. No. Name gender | Religion marks \Rightarrow

selected
gender | marks

Compare student performance

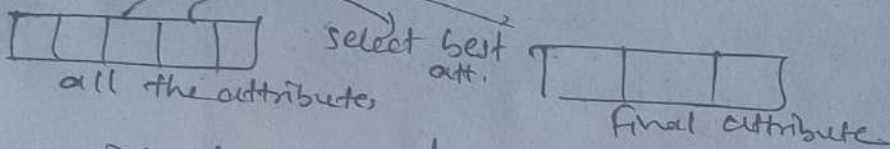
- Exhaustive search tech
(identify each & every combination)
- greedy method
- Best heuristic method

VP → Verb
VP → Verb NP
PP → VP PP
PP → Prep NP

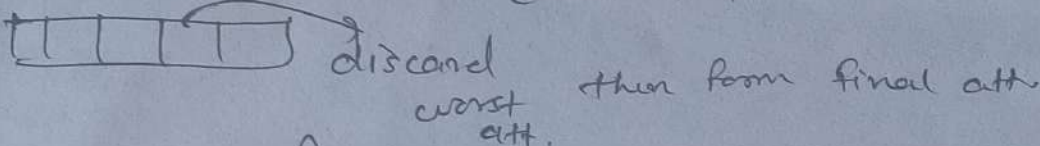
Attribute subset selection
Return original Data

Basic heuristic methods of attribute selection

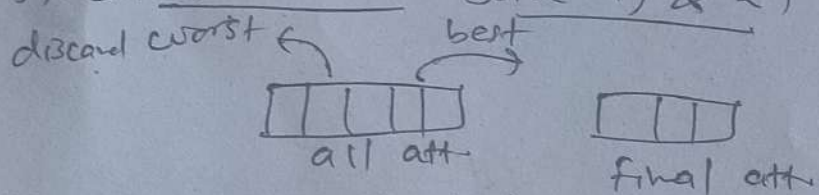
1) Stepwise forward selection



2) Stepwise backward elimination



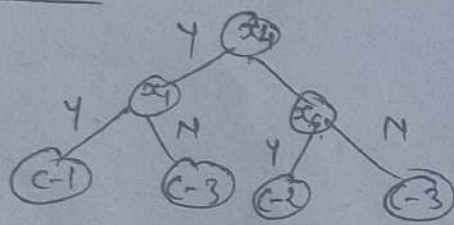
3) Combination of both 1) & 2)



4) decision tree induction : different algo. ⇒ classification

(x_1, \dots, x_6)

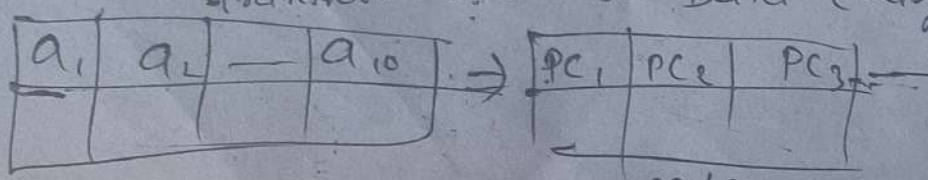
Retain 3 attr
discard all other attr.



Dimensionality Reduction

→ principle component analysis - $N^D \rightarrow n^D$ where $n \leq N$

Transform data so Data change.



2D / 3D

④ Numerosity Reduction

a	b	c

$$y = x_1 + ax_2 + bx_3 + cx_4$$

⇒ equation form through Regression

↳ clustering: grouping data
just store cluster center
instead of storing complete dataset



⑤ data discretization & concept hierarchy generation

(0 - 100)

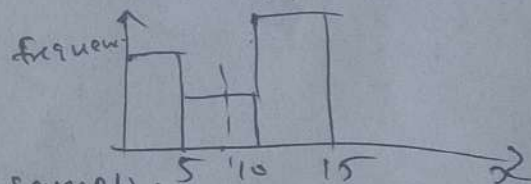
(0-10) (10-30)

break into small parts

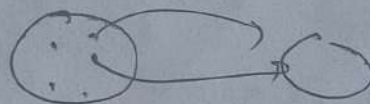
- Best heuristic method

↳ histogram

2, 2, 5, 10, 10, 12, 12, ...



↳ sampling



randomly select few data

10, 20, 30, 50, 60, 90

Young middle age senior

group into higher level

Feature Selection / Attribute Selection Techniques

Why? — Raw data training → model generate
 ↓ why? ↓
 Not relevant features / Noise

Optimal features / relevant features



- Subset of features pass to training
- generate model

① Filter Method: relevant of att.

- Information gain
- chi-square Test
- Correlation coef.

Ex

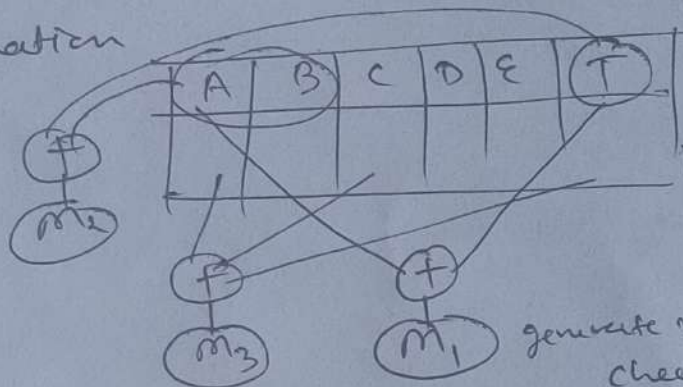
A	B	C	D	E	T	Target attribute
	Roll No				Result	
↓ increase	1 2 3 4				↓ increase	

highly correlated

more important dependency / correlated / relevant att.

② Wrapper Method

- Recursive feature elimination
- Genetic Algo.



take a subset of att. & merge with Target & generate multiple model

- A
- AB
- ABC
- AC

Dis: highly computationally expensive → overfitting problem.

③ Embedded methods } both ② + ③ ⇒ find attribute usefulness

- Decision Tree

less computationally expensive
 overfitting ~~very less~~ does not encounter

☁️ sunny ☁️ cloudy ☔️ raining