

Prove how the given source statement is transformed after each phase of compilation

$x = y + z * 60$

phase 1 - lexical analysis

" 2 - Syntax "

" 3 - semantic " - type info

4 - intermediate code generation

5 - code optimization

$x = y + z * 60$



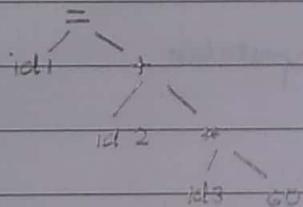
lexical Analysis



$\langle id1, x \rangle \langle op1 \rangle = \rangle \langle id2, y \rangle \langle op2, + \rangle \langle id3, z \rangle \langle op3, * \rangle \langle num, 60 \rangle$



Syntax Analysis



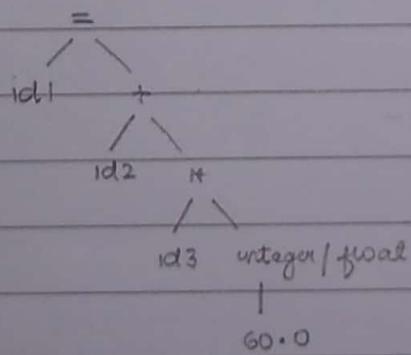
Start $\rightarrow A \text{ start}$

$A \text{ start} \rightarrow id = \text{expr}$

$\text{expr} \rightarrow id + id / id + \text{Num}$



Semantic Analysis



Intermediate code generation



$t_1 = \text{int to float}(60)$

$t_2 = id3 * t_1$

$t_3 = id2 + t_2$

$id1 \leftarrow t_3$



Code optimization

$t_1 = id3 * 60.0$

$id1 = id2 + t_2$



Code generation



MOV F id3, R2
MUL F GO.O, R2
MOVF id2, R1
ADD F R2, R1
MOV F R1, id1

symbol table → has identifiers & their types

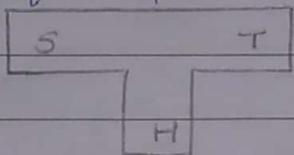


created in 1st phase

but data used in all phases

Bootstrapping and porting

T-diagram of a compiler

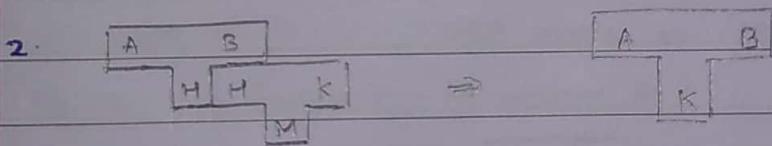
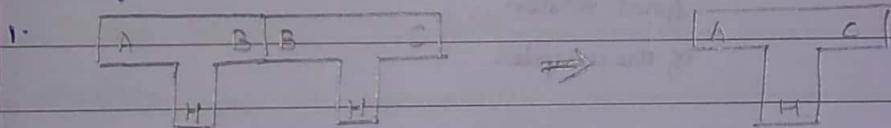


S - source language that it compiles

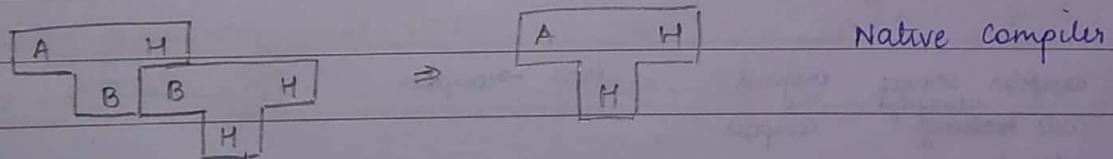
T - target " " " generates code for

H - implementation language that it is written in

T-diagram can be combined in 2 ways

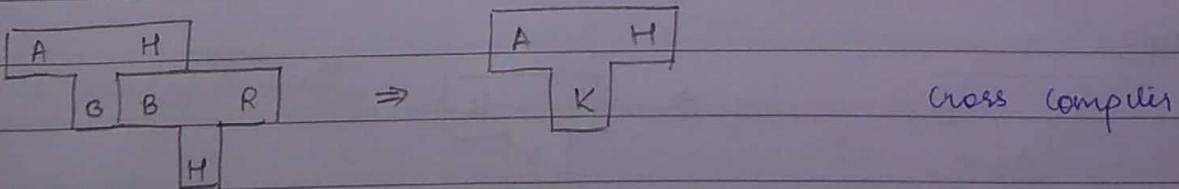


Scenario 1



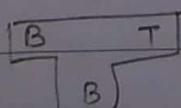
Native compiler

Scenario 2



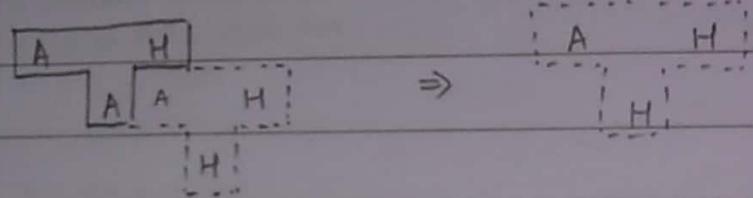
Cross compiler

A compiler with both in the same language that is to compile



circularizing problem - No compiler for the source lang yet exists

Step 1

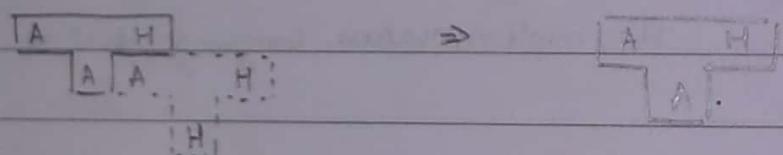


compiler
written in its
own
language

Quick & dirty
compiler
written in
machine
language

Running but
inefficient
compiler

Step 2

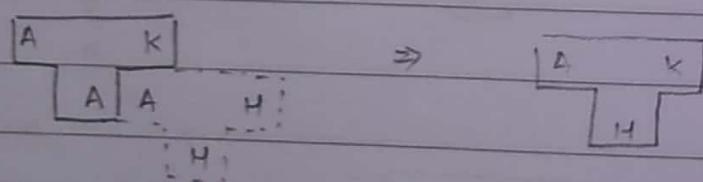


Running but
inefficient
compiler

final version
of the compiler

Porting

Step 1:

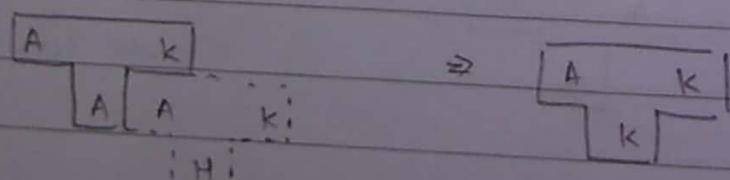


computer service
code ~~retargeted~~
to K

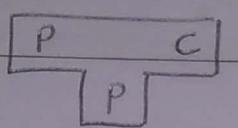
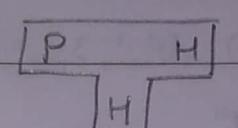
original
computer

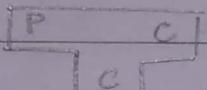
Gross → compiler

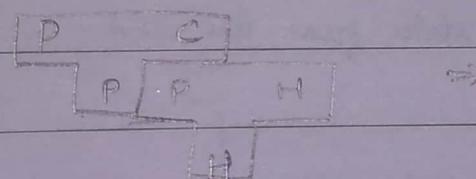
Step 2:

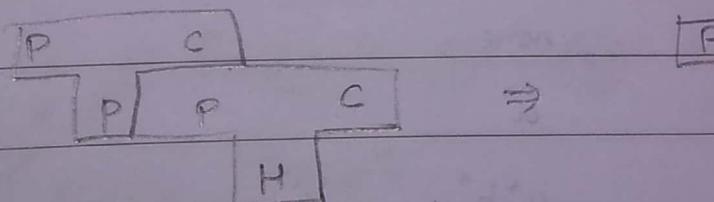


Retargeted
compiler

given :   P - pascal

output 

Step 1 

Step 2 

offunit

UNIT 2 $A \rightarrow A\alpha / \beta$ left recursive $A \rightarrow \alpha A / \beta$ right recursion

elimination of left recursion

(i) $A \rightarrow A\alpha / \beta$

replace with

 $A \rightarrow \beta A'$ $A' \rightarrow \alpha A' / \epsilon$ (ii) $A \rightarrow A\alpha_1 | \dots | A\alpha_m | \beta_1 | \dots | \beta_n$

replace with

 $A \rightarrow \beta_1 A' | \dots | \beta_n A'$ $A' \rightarrow \alpha_1 A' | \dots | \alpha_m A' | \epsilon$ (iii) input : grammar G with no cycles or ϵ production

Output : An equivalent grammar with no left recursion

method :① arrange the NTs in order A_1, A_2, \dots, A_n ② for $i=1$ to n do begin for $j=1$ to $i-1$ do begin replace each production of the form $A_i \rightarrow A_j \gamma$ by the production $A_i \rightarrow S_1 \gamma | \dots | S_k \gamma$ where $A_j \rightarrow S_1 | \dots | S_k$ are all current A_j production and eliminate the immediate left recursion among
 the A_i production

end

egs :

$$\begin{array}{l} A \quad A \alpha \quad \beta \\ \text{if } E \rightarrow E + T / T \\ A \quad A \alpha \quad \beta \\ T \rightarrow T * F / F \end{array}$$
 $F \rightarrow (E) / \text{id} \rightarrow \text{no cycle}$ $G' \quad E \rightarrow TE'$ $E' \rightarrow +TE'/E$ $T \rightarrow FT'$ $T' \rightarrow *FT'/E$ $F \rightarrow (E) / \text{id}$

ii $G: \text{① } S \rightarrow AB$

② $A \rightarrow CB \mid b$

③ $C \rightarrow Sa$

④ $B \rightarrow a$

$S \rightarrow AB \rightarrow \underline{CB} B \rightarrow Sa BB$

$S \rightarrow Sa BB$

\therefore indirect left recursion

Indirect left recursion

$n = 4$

$i = 1 \quad j = 1 \quad A_1 \rightarrow A_1 \quad S \rightarrow S \quad \text{No}$

$i = 2 \quad j = 1 \quad A_2 \rightarrow A_1 \quad A \rightarrow S, \text{ No}$

$i = 3 \quad j = 1 \quad A_3 \rightarrow A_1 \quad C \rightarrow S \quad \text{Yes}$

$C \rightarrow ABA \quad [B \rightarrow AB]$

$j = 2 \quad A_3 \rightarrow A_2 \quad C \rightarrow A \quad \text{Yes}$

$C \rightarrow CBBA \mid bBA$

$A \rightarrow CB$ replace in $C \rightarrow ABA$

Eliminate left recursion from $C \rightarrow CBBA \mid bBC$

$C \rightarrow bBaC'$
 $C \rightarrow BBaC' / E$

$i = 4 \quad j = 1 \text{ to } 3 \quad B \rightarrow S \quad B \rightarrow A \quad C \rightarrow C \quad \text{No}$

grammar becomes

$G^1: \quad S \rightarrow AB$

$A \rightarrow CBb$

$C \rightarrow bBaC$

$C' \rightarrow BBdC' / E$

$B \rightarrow a$

$G : S \Rightarrow (L) / a$

$L \Rightarrow L, S | S$

(i) (a, a)

$S \Rightarrow (L)$
↓
 em

$S \Rightarrow (L)$
↓
 rm

$\Rightarrow (L, S)$
↓
 em

$\Rightarrow (L, S)$
↓
 rm

$\Rightarrow (S, S)$
↓
 em

$\Rightarrow (L, a)$
↓
 rm

$\Rightarrow (a, S)$
↓
 em

$\Rightarrow (S, a)$
↓
 rm

$\Rightarrow (a, a)$
↓
 em

$\Rightarrow (a, a)$
↓
 rm

ii $(a, (a, a))$

left factoring

$$A \rightarrow \alpha B_1 | \alpha B_2 | \gamma$$

Replace with

$$A \rightarrow \alpha A' | \gamma$$

$$A' \rightarrow B_1 | B_2$$

eg G: $A \rightarrow qB | qc$
 $\alpha B_1 | \alpha B_2$

$$A \rightarrow qA' \cancel{qB}$$

$$A' \rightarrow B | C$$

Eq 2
#

$\begin{array}{c} A \\ \alpha \beta_1 \quad \alpha \beta_2 \end{array}$
 $G : \text{state} \rightarrow pQ / pR / \epsilon$

$Q \rightarrow b \text{ stat} | pQQ$

$R \rightarrow p \text{ stat} | bRR$

~~$A \rightarrow pQ$~~

$G' : \text{stat} \rightarrow p \text{ stat}' | \epsilon$

$\text{stat}' \rightarrow Q/R$

$Q \rightarrow b \text{ stat} | pQQ$

$R \rightarrow p \text{ stat} | bRR$

27 observation

Q For full data for attribute age 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33,
33, 35, 35, 35, 35, 35, 36, 40, 45, 46, 52 | Youth
middle age Senior (above 60) |
Sketch egs of each of following sampling techniques SRSWOR, SRSWR, cluster Sampling, Stratified Sampling. Use sample size of 5 and the strata: Youth, middle age, Senior

T ₁	13	T ₁₀	23	T ₁₉	35
T ₂	15	T ₁₁	25	T ₂₀	35
T ₃	16	T ₁₂	25	T ₂₁	35
T ₄	16	T ₁₃	25	T ₂₂	36
T ₅	19	T ₁₄	25	T ₂₃	40
T ₆	20	T ₁₅	30	T ₂₄	45
T ₇	20	T ₁₆	33	T ₂₅	46
T ₈	21	T ₁₇	33	T ₂₆	52
T ₉	22	T ₁₈	35	T ₂₇	70

SRSWOR

Noise - random error or a variance in a measured variable

Eg Attribute values are numeric Eg attribute like price of product
how can we remove this noise thru some data smoothing techniques

① Binning - it smooths a sorted data value by consulting its neighbourhood "value around it". They perform local smooth

- Smoothing by bin means - each value in bin is replaced by the mean value of the bin
- " " " median - replaced by median value of the bin
- " " " boundaries - identify the min & max of each bin

Each value is replaced by the closest boundary value.

② Regression - conform data value to a functⁿ there r 2 types

- linear
- multiple linear regression.

QUESTION

Problem

Suppose that the data for analysis includes the attribute age the age values for the data tuples are (inasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70

(a) Use smoothing by bin means to smooth the above data by using bin depth of 3 & then comment on effect of this technique for given data

- Step 1 sort data

- Step 2: Partition data into $n = \text{depth} = 3$

Bin 1 : 13, 15, 16 Bin 4 : 22, 25, 25 Bin 7 : 35, 35, 35

Bin 2 : 16, 19, 20 Bin 5 : 26, 26, 30 Bin 8 : 36, 40, 45

Bin 3 : 20, 21, 22 Bin 6 : 33, 33, 35 Bin 9 : 46, 52, 70

- Calculate arithmetic mean of each bin

14.667 ; 18.333 ; 21 ; 24 ; 26.667 ; 33.667 ; 35 ; 40.33 ; 56
Bin 1 Bin 2 Bin 3 Bin 4 Bin 5 Bin 6 Bin 7 Bin 8 Bin 9

- Replace each of values in each bin by arithmetic mean

Bin 1 : 14.667, 14.667, 14.667 Bin 4 : 24, 24, 24 Bin 7 : 35, 35, 35

Bin 2 : 18.33, 18.33, 18.33 Bin 5 : 26.667, 26.667, 26.667 Bin 8 : 40.33, 40.33, 40.33

Bin 3 : 21, 21, 21 Bin 6 : 33.667, 33.667, 33.667 Bin 9 : 56, 56, 56

b. use smoothing by bin boundaries using bin depth of 3

Step 1 & 2 same as prev

Bin 1: 13, 16, 16

4: 22, 25, 25

7: 35, 35, 35

2: 16, 20, 20

5: 25, 25, 30

8: 36, 36, 45

3: 20, 20, 22 or 20, 22, 22

6: 33, 33, 35

9: 46, 46, 70

c. What other methods are there for data smoothing

Data Preprocessing

data transformation

data are transformed or consolidated into forms appropriate for mining

Strategies

1. Smoothing - works to remove noise from the data

- Binning
- Regression
- clustering

2. ? ? ?

3. Aggregation - summary or aggregation operations are applied to the data

4. Normalization - attribute data are scaled so as to fall within a smaller range such as -1.0 to 1.0 or 0.0 to 1.0

5. Discretization - raw values of a numeric attribute (e.g. age) are replaced by interval labels (e.g. 0-10, 11-20,...) or conceptual labels (e.g. youth, adult or senior). The labels can be recursively organized into higher level concepts resulting in a concept hierarchy for numeric attribute.

6. Concept hierarchy generation for nominal data - attributes such as street can be generalized to higher level concept like city or country

change from kg to pounds will give diff results
attribute with higher value will have higher weight
Data transformation by Normalization

- this helps avoid dependence on choice of measurements units
 - normalization gives all attributes equal weight
 - useful for classification algorithm & clustering
 - There r many methods for ~~too~~ normalization
- ① Min-max algo

this performs a linear transformation on original data

Suppose \min_A and \max_A are minimum & maximum values of an attribute.

Min-max normalization maps a value v_i of A to v'_i in the range $[\text{new}_{\min}A, \text{new}_{\max}A]$ by computing

$$v'_i = \frac{v_i - \min_A}{\max_A - \min_A} (\text{new}_{\max}A - \text{new}_{\min}A) + \text{new}_{\min}A$$

income has $\min 12000$ Rs

$\max 98000$ Rs

we would like to map income in range $[0.0, 1.0]$

What will 73,600 Rs for income transformed to using min max normalization

$$v'_i = \frac{v_i - \min_A}{\max_A - \min_A}$$

$$= \frac{73600 - 12000}{98000 - 12000} (1.0 - 0.0) + 0$$

$$= \frac{61600}{86000}$$

$$v'_i = 0.716$$

② Z-score normalization (zero mean normalization)

the values for an attribute A are normalized based on the mean (ie average) & the standard deviation of A

The value v_i of A is normalized to v'_i by computing $v'_i = \frac{v_i - \bar{A}}{\sigma_A}$

$$\bar{A} = \text{mean} \quad \sigma_A = \text{standard deviation of attribute A} \quad \sigma_A$$

$$\bar{A} = \frac{1}{n} (V_1 + V_2 + \dots + V_n)$$

$$\sigma_A = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$$

Eg suppose the mean & SD of the values of income are 54000Rs & 16000Rs

suppose what is the value for the income 73600 using z-score normalization?

$$V_i' = \frac{73600 - 54000}{16000} = 1.225$$

③ Normalization by decimal scaling

This normalizes by moving decimal point of values of attribute A the no. of decimal points moved depends on the maximum absolute value of A. A value V_i of A is normalized to V_i' by computing $V_i' = \frac{V_i}{10^j}$ where $j = \text{smallest integer such that}$

$$\max(|V_i'|) < 1$$

Eg suppose that the recorded values of A range from -986 to 917 the max absolute value of A is 986 to normalize by decimal scaling we divide each value by 1000 (ie $j=3$) so that -986 normalizes to -0.986 & 917 normalizes to 0.917

Q Normalize full group of data using (i) min max normalization by setting min = 0 & max = 1 ii) z score normalization

data = 200, 300, 400, 600, 1000

$$i) V_i' = \frac{V_i - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$$

$$200 \quad V_i' = \frac{200 - 200}{1000 - 200} (1.0 - 0.0) + 0.0 = 0$$

$$400 \quad \frac{400 - 200}{1000 - 200} = \frac{200}{800} = 0.25$$

$$300 \quad \frac{300 - 200}{1000 - 200} = \frac{100}{800} = 0.125$$

$$600 \quad \frac{600 - 200}{1000 - 200} = \frac{400}{800} = 0.5$$

$$1000 \quad \frac{1000 - 200}{1000 - 200} = \frac{800}{800} = 1$$

ii 200

$$\bar{A} = \frac{1}{5} (200 + 300 + 400 + 600 + 1000) \\ = 500$$

$$\sigma_A = \sqrt{\frac{1}{5} \sum_{i=1}^5 (x_i - \bar{x})^2} \\ = \sqrt{\frac{1}{5} \left\{ (200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2 + (1000 - 500)^2 \right\}} \\ = \sqrt{\frac{1}{5} \left\{ 90000 + 40000 + 10000 + 100000 + 950000 \right\}} \\ = 282.8$$

$$v_i' = \frac{v_i - \bar{A}}{\sigma_A}$$

200	$v_1' = -1.05$
300	$= -0.17$
400	$= 0.35$
600	$= 0.35$
1000	$= 1.72$

Discretization and Binarization

transform a continuous attribute into a categorical attribute (discretization) and both continuous & discrete attributes need to be transformed into 1 or more binary attributes (binarization)

Binarization

A simple technique of Binarization of a category will be as follows if there are m categorical values then uniquely assign each original value to the integer in the interval $[0, m-1]$ then convert each of these m integers to a binary no. Since $n = \lceil \log_2 m \rceil$ binary digits are required to represent the integers the binary nos are represented using n binary attributes

e.g. categorical variable with 5 values {awful, poor, ok, good, great}

\therefore there r 5 values \rightarrow would require 3 binary variables x_1, x_2, x_3 & conversion is shown in table

attribute type

disimilarity

similarity

Nonnumerical

$$d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$$

$$s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$$

product measure
ordinal

$$d = |x - y| / (n - 1) \quad (\text{values})$$

$$s = 1 - d$$

mapped to integers 0 to $n-1$,
where $n = \text{no. of values}$)

③ eg. ht, wt,

Interval ratio

$$d = |x - y|$$

$$s = -d$$

$$s = \frac{1}{1+d}$$

$$s = e^{-d}$$

last year 1 m. 72

5 kg more than

than year

75

$$s = \frac{1 - d - \min_d}{\max_d - \min_d}$$

Euclidean Distance

distance b/w 2 objects with n no. of attributes

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

$$P_1 = (0, 2)$$

$$P_2 = (2, 0)$$

$$P_3 = (3, 1)$$

$$P_4 = 5.71$$

$$P_1 P_3 = \sqrt{(-3)^2 + (1)^2} = \sqrt{10}$$

$$P_1 P_4 = \sqrt{25 + 1} = \sqrt{26} = 5.099$$

$$P_2 P_1 = \sqrt{2^2 + (-2)^2} = \sqrt{8}$$

	P_1	P_2	P_3	P_4	
P_1	0	2.828	3.162	5.099	$P_2 P_4 = \sqrt{9 + 1} = 3.162$
P_2	2.828	0	1.414	3.162	$P_3 P_4 = \sqrt{2^2 + 2^2} = \sqrt{8} = 2$
P_3	3.162	1.414	0	2	
P_4	5.099	3.162	2	0	

$$d(p_1, p_2) = \sqrt{(0 - 2)^2 + (2 - 0)^2} \approx \sqrt{2^2 + 2^2} = \sqrt{8}$$

$$= 2.828$$

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

$$d(p_1, p_2) = d(p_2, p_1) \quad d(p_1, p_1) = 0$$

Minkowski Distance - generalization of euclidean

$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^n \right)^{1/n}$$

e.g. $P_1(0, 2)$ $P_2(2, 0)$ $P_3(3, 1)$ $P_4(5, 1)$

$L_1 \quad n = 1$

	P_1	P_2	P_3	P_4
P_1	0	4	4	6
P_2	4	0	2	4
P_3	4	2	0	2
P_4	6	4	2	0

$L_2 \quad n = 2$

	P_1	P_2	P_3	P_4
P_1	0	2.828	3.162	5.099
P_2	2.828	0	1.414	3.162
P_3	3.162	1.414	0	2
P_4	5.099	3.162	2	0

$L_\infty \quad n = \infty$

	P_1	P_2	P_3	P_4
P_1	0	2	3	5
P_2	2	0	1	3
P_3	3	1	0	2
P_4	5	3	2	0

Common properties of distance

- ① $d(x, y) \geq 0$ $\forall x$ and y and $d(x, y) = 0$ iff $x = y$
 - ② $d(x, y) = d(y, x)$ $\forall x$ and y (symmetry)
 - ③ $d(x, z) \leq d(x, y) + d(y, z)$ \forall points x, y and z (triangle inequality)
- where $d(x, y)$ = distance (dissimilarity) between points (data objects) x and y

Similarity betⁿ binary vector

patient table		fever	cough	test 1	test 2	test 3	test 4
name	gender						
Jack	M	Y	Y	P	N	P	N
John	M	Y	N	N	N	N	P
Mary	F	Y	Y	P	N	P	P
?	?	?	?	?	?	?	?

Simple Matching and Jaccard coefficients

$$\text{SMC} = \frac{\text{no. of matches}}{\text{no. of attributes}}$$

$$= \frac{(f_{11} + f_{00})}{(f_{01} + f_{10} + f_{11} + f_{00})}$$

f_{00} = no. of attributes when $x=0, y=1$

$f_{10} =$

$x=1, y=0$

$f_{01} =$

$x=0, y=1$

$f_{11} =$

$x=1, y=1$

$$J = \frac{\text{no. of 11 matches}}{\text{no. of non 0 attributes}}$$

$$= \frac{f_{11}}{f_{01} + f_{10} + f_{11} + f_{00}}$$

- when x have asymmetric J plays an imp role
- when x have symmetric SMC plays an imp role

Jaccard coefficient (Jack, John)

$$\frac{1}{1+3+1} = \frac{1}{5} = 0.2$$

assume yes = 1
no = 0

+ve
 $P=1$
-ve
 $N=0$

similarity b/w doc1 & doc2 (cosine of vectors) is $\theta \approx 90$ far from each other. $\theta \approx 0 \rightarrow$ more similar.

Cosine Similarity

Document Vector

document	team	coach	hockey	baseball	baseball	soccer	penalty	score	win	loss	season
doc 1	5	0	3	0	2	0	0	2	0	0	0
doc 2	3	0	2	0	1	1	0	1	0	1	1
doc 3	0	7	0	2	1	0	0	3	0	0	0
doc 4	0	1	0	0	1	2	2	0	3	0	0

Term frequency vector or $(5, 0, 3, 0, 2, 0, 0, 2, 0)$ is the term freq vector of doc 1

Let $x \& y$ be 2 vectors for comparison

similarity is given by

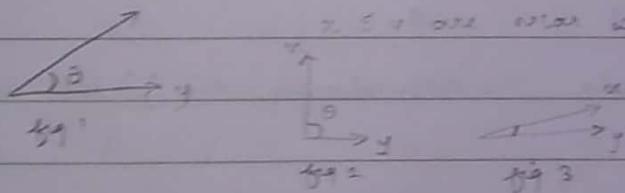
$$\text{Sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

where $\|x\|$ = Euclidean norm of vector $x = (x_1, x_2, \dots, x_n)$ defined as

$\sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$ conceptually as the length of vector

similarly $\|y\|$ is the Euclidean norm of vector y

compute the cosine of the angle between $x \& y$



smaller angles: they are more similar

similarity bet doc1 and doc 2

$$x = (5, 0, 3, 0, 2, 0, 0, 2, 0) = \text{term freq of doc 1}$$

$$y = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1) = " " " " " 2$$

$$\begin{aligned} x \cdot y &= 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 0 + 2 \times 1 + 0 \times 0 + 0 \times 1 \\ &= 15 + 6 + 2 + 2 = 25 \end{aligned}$$

$$\begin{aligned} \|x\| &= \sqrt{25 + 0 + 9 + 0 + 4 + 0 + 0 + 4 + 0 + 0} \\ &= \sqrt{25 + 9 + 4 + 4} = \sqrt{42} = 6.48 \end{aligned}$$

$$\begin{aligned} \|y\| &= \sqrt{9 + 0 + 4 + 0 + 1 + 1 + 0 + 1 + 0 + 1} \\ &= \sqrt{9 + 4 + 4} = \sqrt{17} = 4.12 \end{aligned}$$

$$\text{Sum}(x, y) = \frac{25}{6.48 \times 4.12} = \frac{25}{26.697} = 0.94 \quad \begin{matrix} \text{very close to 1} \\ \therefore \text{similar} \end{matrix}$$

$\text{Sum}(\text{doc 1}, \text{doc 4})$

$$x = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$y = (0, 1, 0, 0, 1, 2, 2, 0, 3, 0)$$

$$x \cdot y = 2$$

$$\|x\| = 6.48$$

$$\|y\| = \sqrt{1+1+4+4+9} = \sqrt{19} = 4.358$$

$$\text{Sum}(x, y) = \frac{2}{6.48 \times 4.358} = \frac{2}{28.239} = 0.0708$$

Dissimilarity for attributes of mixed type

Object identifier	Test 1 (nominal)	Test 2 (ordinal)	Test 3 (numeric)
1	Code A	Excellent	45
2	Code B	Fair	22
3	Code C	Good	64
4	Code A	Excellent	28

In such situations we use the formula: The dissimilarity $d(i, j)$ between objects $i \neq j$ is defined as

$$d(i, j) = \frac{\sum_{f=1}^P s_{ij}^{(f)} \cdot d_{ij}^{(f)}}{\sum_{f=1}^P s_{ij}^{(f)}}$$

where we $s_{ij}^{(f)} = 0$ if either:

x_{if} or x_{jf} is missing (ie there is no measurement of attribute f for object i or object j)

OR $x_{if} = x_{jf} = 0$ the attribute f is asymmetric binary otherwise $s_{if}^{(f)} = 1$

the contribution of attribute f to the dissimilarity between $i \neq j$ ie $d_{ij}^{(f)}$ is computed depending on its type

① if f is numeric $d_{ij}^{(f)}$ is computed using $d_{ij}^{(f)} = \frac{|x_{if} - x_{jf}|}{\max x_{hf} - \min x_{hf}}$

where h runs over all non missing objects for attribute f

② if f is nominal or binary $d_{ij}^{(f)}$ $\begin{cases} 0 & \text{if } x_{if} = x_{jf} \\ 1 & \text{otherwise} \end{cases}$

③ if $f =$ ~~nominal~~ ^{ordinal} compute the ranks r_{if} & $z_{if} = \frac{r_{if} - 1}{M_f - 1}$ treat z_{if} as numeric

No. of values
 $3 - 1 = 2$ for our ex

	1	2	3	4	dissimilarity matrix for test 3	
1	0					
2	0.55	0				
3	0.45	1	0			
4	0.40	0.14	0.86	0		

dissimilarity for test 2

	1	2	3	4		
1	0				fair - 0	
2	1	0			good - 1	
3	0.5	0.5	0		excellent - 2	
4	0	1	0.5	0	<u>2 - 0</u>	<u>2 - 1</u>

$2 - 0 = 2$ $2 - 1 = 1$ $1 - 0 = 1$ $0 - 0 = 0$

dissimilarity for test 1

	1	2	3	4
1	0			
2	1	0		
3	1	1	0	
4	0	1	1	0

Code B, A

Code A, Code B

dissimilarity	1	2	3	4
1	0			
2	<u>0.85</u>	0		
3	0.65	0.83	0	
4	0.13	0.71	0.79	0

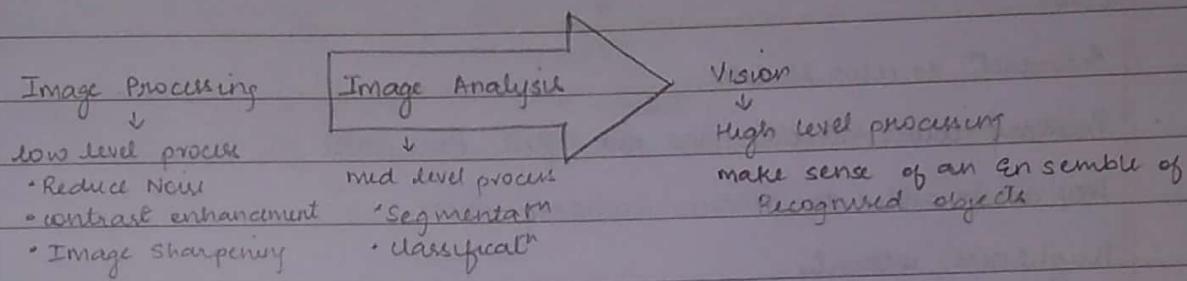
$$d(4,1) = 0.4/3$$

$$d(4,2) 0.114 + 2 = 0.21$$

Objects 1 & 4 r most similar

" 2 & 1 r " dissimilar

Image processing



Application

- Medical Images
- Face Recognition
- Finger print Detection
- Pattern Recognition
- Feature extraction

Definitions :-

- Image - An image is a 2-D signal analog or digital that contains intensity or colour information arranged along the x and y spatial axis
general aspect
of with respect to
It is a 2D function $f(x, y)$ where x, y are spatial coordinates & f is amplitude of a function also called as the intensity or grey level or the color of the image at that point

0	0	0	0	0	0	0
1	1	0	1	1	1	1
1	1	1	0	1	1	1
1	1	1	0	1	1	1
1	1	1	0	1	1	1

0 - represents black
1 - " white
these are called a monochrome scheme

- Amplitude - measure of how far and in what direction a variable differs from 0
- Image Processing - IP is a form of signal processing for which the input is an image such as photograph or frames or video and opp can be an image or set of characteristics or parameters related to the image
- In image processing we do need it \rightarrow its a method of performing operation on an image to get an enhanced image or to extract useful info from it
- digital image processing is a subset of digital signal processing
- it allows wider range of algo to be applied to the ip data which can avoid problems

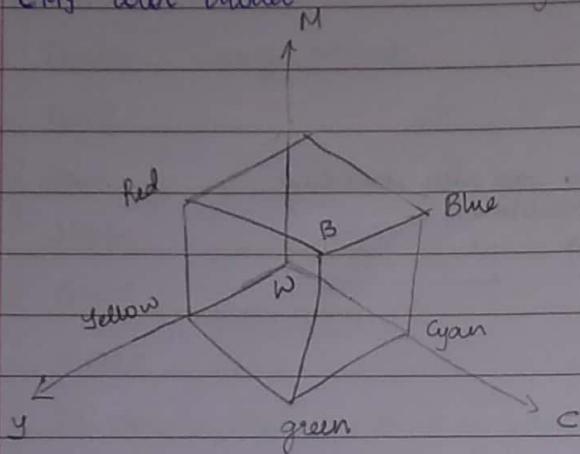
of noise & distortion (no image sharpening)

Color

Introduction to color Models

- Primary colors - absorbs and also reflect other colors
they can be added to produce secondary colors
- Brightness - intensity
- ~~Hue~~ Hue - is a dominant wavelength (color)
- Saturation - amt of white light mixed with the hue
- Chromaticity - depends on a combination of hue and saturation
- Luminosity - ~~maximum~~ ^{measure of} light energy emitted (More brighter - more energy - more luminosity)
- Luminance - measure of luminous intensity per unit area of light off axis

CMY Color Model C- cyan M- magenta Y- yellow



- this model is used for describing the color o/p to hardcopy device
- colors are specified by what is subtracted from white color

$$\text{Green} + \text{Blue} = \text{Cyan}$$

when white light is reflected from cyan color mean the reflected light doesn't have the red component ie red light is absorbed or subtracted

$$\begin{bmatrix} C \\ M \\ Y \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} R \\ G \\ B \end{bmatrix}$$

$$\begin{bmatrix} R \\ G \\ B \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} C \\ M \\ Y \end{bmatrix}$$

$$I_{CMY} = (I_C \quad I_M \quad I_Y)$$

$$I_C = 1 - \frac{f_R(x,y)}{f_R(x,y) + f_G(x,y) + f_B(x,y)}$$

$$I_M = 1 - \frac{f_G(x,y)}{f_R(x,y) + f_G(x,y) + f_B(x,y)}$$

$$I_Y = 1 - \frac{f_B(x,y)}{f_R(x,y) + f_G(x,y) + f_B(x,y)}$$

Convert RGB values to CMY model

RGB values (29, 98, 128)

$$\begin{array}{r} 2 \\ 128 \\ 98 \\ \hline 29 \\ 285 \\ \hline 229 \end{array}$$

$$\cancel{dc} = 1 - \frac{29}{255} = \frac{255 - 29}{255} = \frac{226}{255} = 0.88$$

$$\frac{29}{255} = 0.113$$

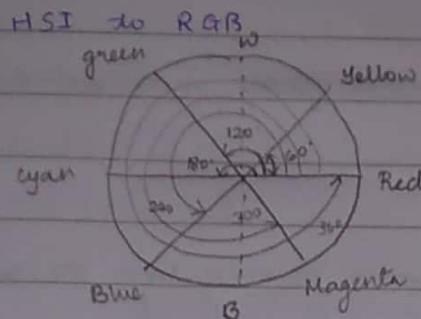
$$\frac{98}{255} = 0.384$$

$$\frac{128}{255} = 0.5019$$

$$\begin{aligned} Cyan &= 1 - R \\ &= 1 - 0.113 \\ &= 0.887 \end{aligned}$$

$$\begin{aligned} Magenta &= 1 - G \\ &= 1 - 0.384 \\ &= 0.616 \end{aligned}$$

$$\begin{aligned} Yellow &= 1 - B \\ &= 1 - 0.501 \\ &= 0.499 \end{aligned}$$



angles $r = \text{placed}$
 Red = 0°
 Yellow = 60° green = 120° Cyan = 180° Blue = 240°
 Magenta = 300° Red = 360°

Note where $H = ?$

$H = 0$ at $I = 0$ is red color

If $0 < H < 120^\circ$ Blue not present $\therefore I = IS$

$$R = I + 2IS \quad *$$

$$R = I + IS \times \frac{\cos(H)}{\cos(60 - H)}$$

$$G = I - IS$$

$$G = I + IS \times [1 - \frac{\cos(H)}{\cos(60 - H)}]$$

$$B = I - IS$$

$$B = I - IS$$

$H = 120^\circ$ is green

If $120 < H < 240^\circ$ Red and magenta

$$R = I - IS$$

$$R = I - IS$$

$$G = I + 2IS \quad *$$

$$G = I + IS \times [\cos(H - 120) / \cos(180 - H)]$$

$$B = I - IS$$

$$B = I + IS \times [1 - \cos(H - 120) / \cos(180 - H)]$$

$H = 240^\circ$

If $240 < H < 360^\circ$ Green and cyan

$$R = I - IS$$

$$R = I + IS \times [1 - \cos(H - 240) / \cos(300 - H)]$$

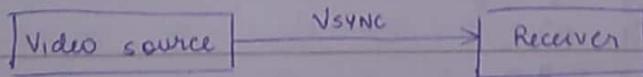
$$G = I - IS$$

$$G = I - IS$$

$$B = I + 2IS \quad *$$

$$B = I + IS \times [\cos(H - 240) / \cos(300 - H)]$$

Video



① VSYNC generation

② 1st line image is scanned

③ level synchronisation signal

- It's a sequence of continuous pictures and each picture in the sequence is called as a frame

• 25 frames or more per second ~~so that each frame must be displayed~~

→ term used to compare 2 images.

iii a process of moving a filter mask over the image in computing the sum of the product at each location

$$\text{let } I = \{0, 0, 1, 0, 0\} \quad \text{mask} = \{3, 2, 8\}$$

$$\begin{matrix} 0 & 0 & 1 & 0 & 0 \end{matrix} \Rightarrow \begin{matrix} 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{matrix}$$

$$3 \ 2 \ 8$$

$$3 \ 8 \ 2$$

$$3 \ 2 \ 8$$

$$3 \ 2 \ 8$$

(Step i) a 0 padding process

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & & & & & & & 0 \end{matrix}$$

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & & & & & & & 0 \end{matrix}$$

input in order of mask

ii initial position step

$$a_0 + 2 \cdot 0 + 8 \cdot 0 = 0$$

iii right shift by 1 bit position

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & & & & & & & 0 \end{matrix}$$

iv

right shift by 2 bit

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 8 & & & & & \end{matrix}$$

v right shift by 3 bit

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 8 & 2 & & & & \end{matrix}$$

vi right shift by 4 bit

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 8 & 2 & 3 & & & \end{matrix}$$

vii right shift by 5 bit

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 8 & 2 & 3 & 0 & & \end{matrix}$$

viii right shift by 6 bit

$$\begin{matrix} 3 & 2 & 8 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 8 & 2 & 3 & 0 & 0 & & \end{matrix}$$

convolution

$$I = \{0, 0, 1, 0, 0\}$$

$$K = \{3, 2, 8\}$$

step i 0 padding process

Rotate Mask = 180°

$$K = \{8, 2, 3\}$$

ii initial slip

$$\begin{matrix} & 8 & 2 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ \bullet & & & & & & & & \end{matrix}$$

iii right shift by 1 position

$$\begin{matrix} & 8 & 2 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ \bullet & \bullet & & & & & & & \end{matrix}$$

iv right shift by 2 position

$$\begin{matrix} & & 8 & 2 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ \bullet & \bullet & 3 & & & & & & \end{matrix}$$

v

$$\begin{matrix} & & & 8 & 2 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 2 & & & & & \end{matrix}$$

vi

$$\begin{matrix} & & & & 8 & 2 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 2 & 8 & & & & \end{matrix}$$

vii

$$\begin{matrix} & & & & & 8 & 2 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 2 & 8 & 0 & & & \end{matrix}$$

viii

$$\begin{matrix} & & & & & & 8 & 2 & 3 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 2 & 8 & 0 & 0 & 0 & 0 \end{matrix}$$

2 W

$$I = \begin{Bmatrix} 3 & 3 \\ 3 & 3 \end{Bmatrix} \text{ be an image } \& K = \begin{Bmatrix} 1 & 2 \\ 3 & 4 \end{Bmatrix}$$

be a kernel {mask} Perform convolution & correlation

convolution

① Rotate by 180°

$$\textcircled{1} \quad \begin{Bmatrix} 1 & 2 \\ 3 & 4 \end{Bmatrix} \text{ rotate by } 180^\circ \xrightarrow{\textcircled{2}} \begin{Bmatrix} 3 & 4 \\ 1 & 2 \end{Bmatrix} \Rightarrow \textcircled{3} \quad \begin{Bmatrix} 4 & 3 \\ 2 & 1 \end{Bmatrix}$$

(a)

$$\textcircled{4} \quad \begin{Bmatrix} 3 & 3 \\ 3 & 3 \end{Bmatrix} \quad \boxed{\begin{matrix} 4 & 3 \\ 2 & 1 \end{matrix}} \text{ kernel}$$

$$\textcircled{5} \quad \begin{matrix} 0^4 & 0^3 \\ 0^2 & 3^1 \end{matrix} \quad 0 \quad 0 \quad \textcircled{6} \quad (4 \times 0) + (3 \times 0) + (2 \times 0) + (3 \times 1) = 3 \quad \text{o/p written on top set of mask} \\ 0 \quad 3 \quad 3 \quad 0 \quad \textcircled{7} \quad \begin{matrix} 3 & 0 & 0 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

$$\textcircled{8} \quad \text{(b)} \quad \begin{matrix} 3 & \boxed{0^4 0^3} & 0 \\ 0 & \boxed{3^2 3^1} & 0 \\ 0 & 3 & 3 \quad 0 \\ 0 & 0 & 0 \quad 0 \end{matrix} \quad \textcircled{9} \quad 4 \times 0 + 0 \times 3 + 3 \times 2 + 3 \times 1 = 9 \\ \textcircled{10} \quad \begin{matrix} 3 & 9 & \boxed{0^4 0^3} \\ 0 & 3 & \boxed{3^2 0^1} \\ 0 & 3 & 3 \quad 0 \\ 0 & 0 & 0 \quad 0 \end{matrix} \quad \textcircled{11} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 0^4 3^3 & 3 & 0 & 0 \\ 0^2 3^1 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \quad \textcircled{12} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 0^4 3^3 & 3 & 0 & 0 \\ 0^2 3^1 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \quad \textcircled{13} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

$$\textcircled{14} \quad \text{(d)} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 0^4 3^3 & 3 & 0 & 0 \\ 0^2 3^1 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

$$\textcircled{15} \quad \text{(e)} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 1^2 2 & \boxed{3^4 3^3} & 0 & 0 \\ 0 & \boxed{3^2 3^1} & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \quad \begin{matrix} 3 \times 4 + 3 \times 3 \\ 3 \times 2 \times 3 \times 1 \\ 12 + 9 + 6 + 3 \\ 30 \end{matrix} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 12 & 30 & 3 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

$$\textcircled{16} \quad \text{(f)} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 12 & 30 & \boxed{3^4 0^1} & 0 \\ 0 & 3 & \boxed{3^2 0^1} & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \quad \begin{matrix} 3 \times 4 + 0 \times 3 \\ 3 \times 2 + 0 \times 1 \\ 18 \end{matrix} \quad \begin{matrix} 3 & 9 & 6 & 0 \\ 12 & 30 & 18 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

(g)

3	9	6	0
12	30	18	0
0 ⁴	3 ³	3	0
0 ²	0 ¹	0	0

$0 \times 4 + 3 \times 3 + 0 \times 2$
 $+ 0 \times 1$

3	9	6	0
12	30	18	0
9	3	3	0
0	0	0	0

(h)

3	9	6	0
12	30	18	0
9	21	3 ⁴ 3 ³	0
0	0 ²	0 ¹	0

$3 \times 4 + 3 \times 3$
 $0 \times 2 + 0 \times 1$

3	9	6	0
12	30	18	0
9	21	3	0
0	0	0	0

(i)

3	9	6	0
12	30	18	0
9	21	3 ⁴ 0 ³	0
0	0	0 ² 0 ¹	0

$3 \times 4 + 3 \times 0$
 $+ 0 \times 2 + 0 \times 1$

3	9	6	0
12	30	18	0
9	21	12	0
0	0	0	0

What is the highest value? 30 have placed mark on the area which contains no 1

Correlation

$$\begin{Bmatrix} 1 & 2 \\ 3 & 4 \end{Bmatrix}$$

0 ¹ 0 ²	0	0
0 ³ 3 ⁴	3	0
0	3	3 0
0	0	0 0

$0 \times 1 + 0 \times 2$
 $+ 0 \times 3 + 3 \times 4$
 $= 12$

12	0 ¹ 0 ²	0
0	3 ³ 4 ³	0
0	3	3 0
0	0	0 0

$0 \times 1 + 0 \times 2$
 $+ 3 \times 3 + 3 \times 4$
 $= 9 + 12$
 $= 21$

12	21	0 ¹ 0 ²
0	3 ³ 0 ⁴	0
0	3	3 0
0	0	0 0

0 ¹ 0 ²	0	0
0 ³ 3 ⁴	3	0
0	3	3 0
0	0	0 0

$0 \times 1 + 0 \times 2$
 $+ 3 \times 3 + 0 \times 4$
 $= 0 + 9 + 0$
 $= 9$

12	21	9	0
0 ¹ 2 ³	3	0	
0 ³ 4 ³	3	0	
0	0	0	

$0 \times 1 + 2 \times 3 +$
 $0 \times 3 + 4 \times 3$
 $= 6 + 12$
 $= 18$

12	21	9	0
0	3 ¹ 2 ³	0	
0	3 ³ 4 ³	0	
0	0	0	

$3 \times 1 + 2 \times 3$
 $+ 3 \times 3 + 4 \times 3$
 $= 3 + 6 + 9 + 12$
 $= 30$

12	21	9	0
18	30	3 ¹ 2 ⁰	0
0	3	3 ³ 4 ⁰	0
0	0	0	0

$3 \times 1 + 3 \times 3$
 $= 3 + 9$
 $= 12$

12	21	9	0
18	30	3 ¹ 2 ³	0
0	3	3 ³ 4 ³	0
0	0	0	0

$3 \times 2 +$
 $= 6$

12	21	9	0
18	30	3 ⁰ 4 ¹ 2 ²	0
6	3 ⁴ 2 ³	0	0
0	0 ³ 4 ⁰	0	0

$3 \times 1 + 2 \times 3$
 $= 3 + 6$
 $= 9$

12	21	9	0
18	30	12	0
6	9	3 ¹ 2 ⁰	0
0	0	0 ³ 4 ⁰	0

$3 \times 1 + 0$
 $= 3$

12	21	9	0
18	30	12	0
6	9	3	0
0	0	0	0

Convolution Eq

$$x(m) \cdot h(g-m)$$

length of 1st sequence = l_1

" " 2nd " = l_2

" " output " = $l_1 + l_2 - 1$

First sequence starts at $n = n_1$

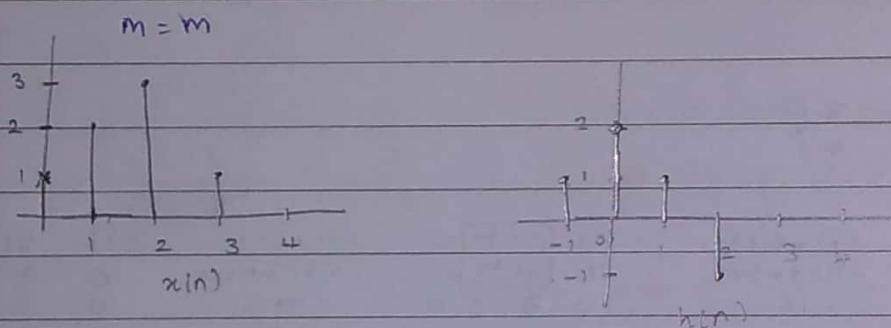
Second sequence starts at $n = n_2$

Output sequence starts at $n = n_1 + n_2$

The " " ends at $(n_1 + n_2) + (l_1 + l_2 - 2)$

Q. $x(n) = \{ \begin{smallmatrix} 0 & 1 & 2 & 3 \\ 1 & 2 & 3 & 1 \end{smallmatrix} \}$ $h(n) = \{ \begin{smallmatrix} -1 & 0 & 1 & 2 \\ 1 & 2 & 1 & -1 \end{smallmatrix} \}$

①



$$l_1 = 4 \quad l_2 = 4$$

$$\text{result length} = l_1 + l_2 - 1 = 4 + 4 - 1 = 7$$

$$n_1 = 0 \quad n_2 = -1 \quad n = n_1 + n_2 = 0 - 1 = -1$$

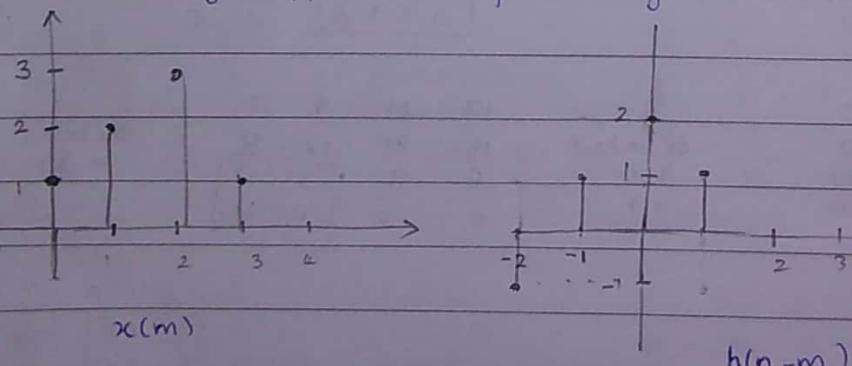
off end at $(n_1 + n_2) + (l_1 + l_2 - 2)$

$$(0 - 1) + (4 + 4 - 2) = 5 \quad \textcircled{2} \quad x(n) = x(m)$$

③

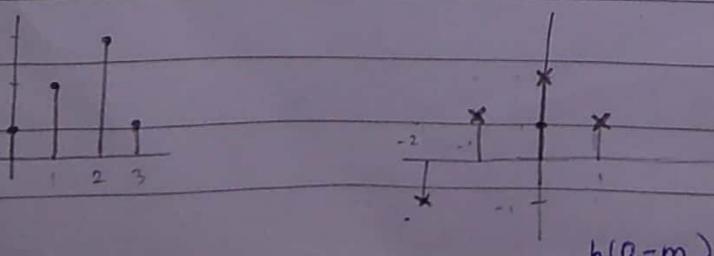
$m = -m$ change happens \Rightarrow only in h signal

④

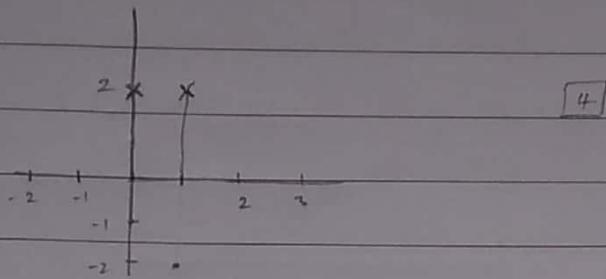


④

$$g = 0$$

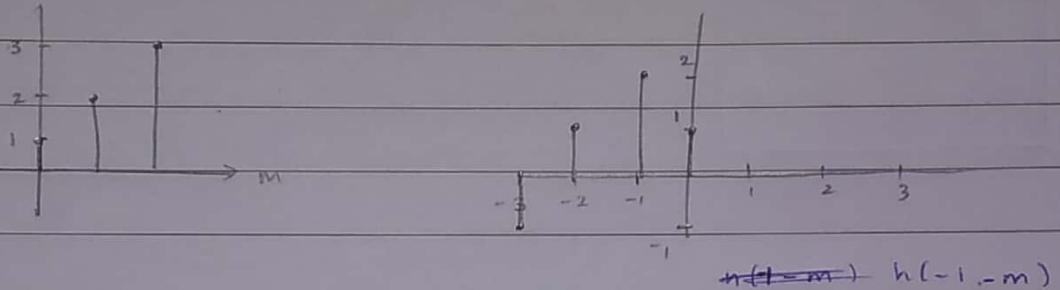


$$x(m) \cdot h(0-m)$$

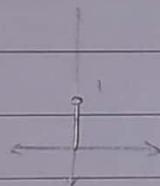


(5)

$$q = -1$$



$$x(m) \cdot h(-1-m)$$



✓ final answer $\{1, 4, 8, 8, 3, -2, -1\}$

