

Mod 5 - Unsupervised learning

classmate

Date _____

Page _____

unsupervised learning

Anomaly detection → clustering → Association
→ association mining → rule mining.
→ univariate based → bivariate based
→ multi-variate → positive test
→ negative test
→ outlier detection

1) Anomaly detection:-

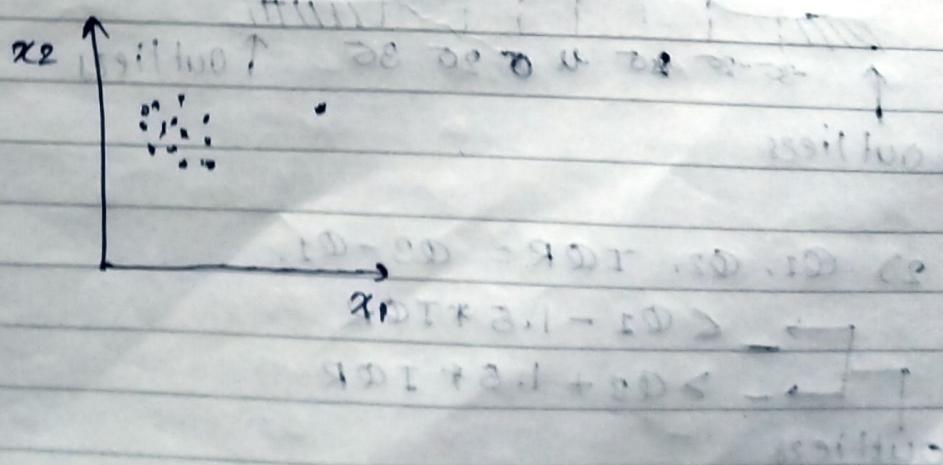
- 1) distance based
- 2) density based
- 3) Model based.

• Outliers:-

- data points that significantly differ from data points in dataset.
- These data points can distort statistical analyses and ml models, leading to inaccurate predictions or conclusions.

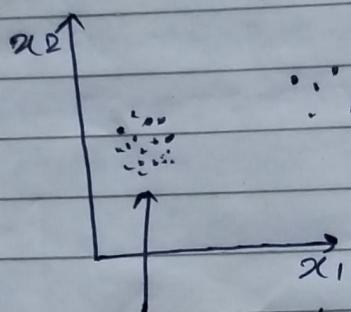
1) Distance based:-

- this method identifies outliers based on the distance of data points from others in dataset.
- These methods assume that outliers are located far away from the majority of data points.



2) Density based

- identifies outliers by assessing the density of data points within the dataset.
- These methods are based on the assumption that outliers are data points that lie in low-density regions compared to majority of data.



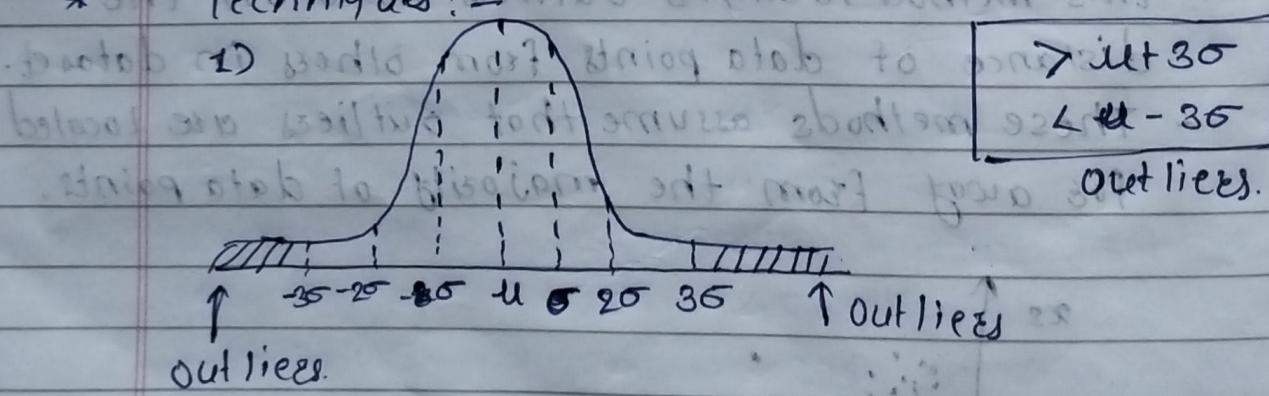
most left side is high dense. rest is low dense -

rest of the data is low dense

3) Model based

- higher ml models used to determine outliers. ex. Decision tree (iforest), Regression.

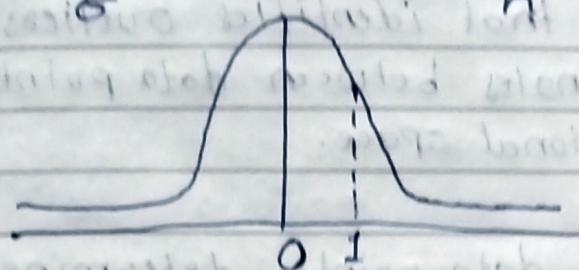
* Out Techniques:



$$2) Q_1, Q_3, IQR = Q_3 - Q_1$$

$$\begin{cases} \text{outlier} & < Q_1 - 1.5 * IQR \\ \text{outlier} & > Q_3 + 1.5 * IQR \end{cases}$$

$$3) z_i = \frac{x_i - \mu_{\text{all data}}}{\sigma} \sim N(0, 1)$$



2173

$$z_1 < -3$$

outlines.

4) Median absolute deviation (MAD): -

$$MAD = \sum_{i=1}^n |x_i - m| \quad \text{and} \quad m = \text{median}$$

Given, data: - Intertwined basis tables

$$f = 5, 10, 2, 3, 25, 45, 100$$

Sqet, 2, 3, 5, 10, 25, 45, 100
M

$$\therefore \text{MAD} = (5-10) + (10-10) + (2-10) + (3-10) + (25-10) \\ + (45-10) + (100-10)$$

$$\left. \begin{array}{l} \rightarrow M - C \cdot MAD \\ \downarrow \rightarrow M + C \cdot MAD \end{array} \right\} \text{outliers.}$$

5) K-NN: (kth nearest neighbour) usually k = 3

* Techniques to find NN.

1) L1 | Manhattan | city block.

$$\rightarrow L^1 = |x_1 - x_2| + |y_1 - y_2| \dots$$

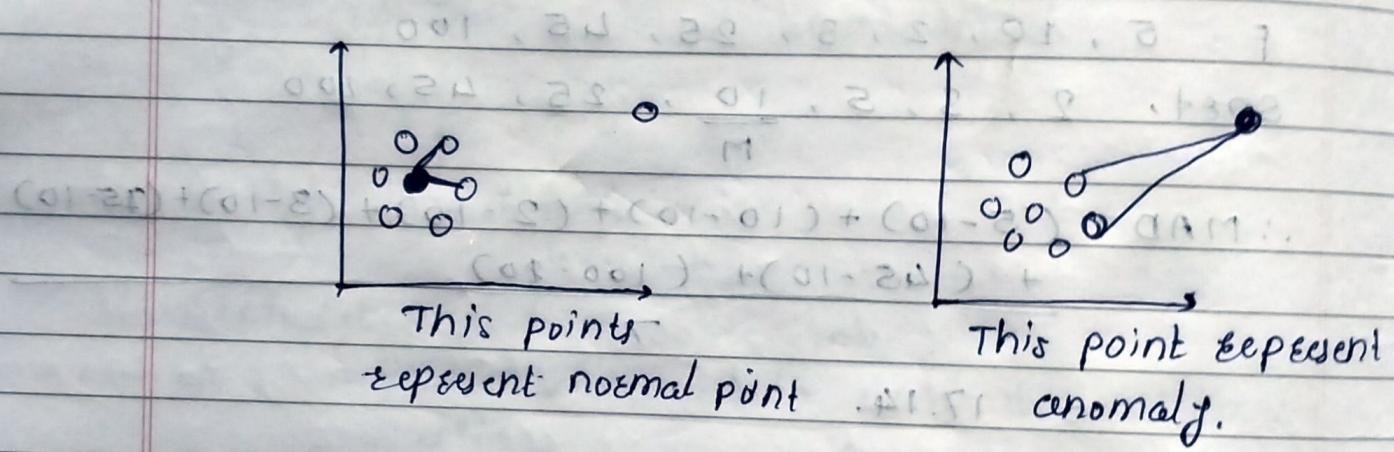
2) L2 / distance / Euclidian :-

$$\rightarrow L_2 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

- * Angle based anomaly detection:
 - method that identifies outliers by measuring the angles between data points in high-dimensional space.

steps :-

- 1) For each data point, determine the angle it makes with all pairs of other data points
- 2) Calculate (the variance) of this angle.
- 3) Points for which the variance is below predetermined threshold are anomalies.



- * Example on k-NN :-
- Determine the outliers for threshold, for distance kth neighbour as 22. Use L1 measure. $k=8$.

ID	F1	F2	F3
1	12	22	48
2	11	123	48
3	3	21	47
4	4	22	46
5	8	24	49
6	12	19	51
7	11	38	55
8	10	19	56
9	13	20	49
10	14	20	49

$$L_1 = |x_2 - x_1| + |y_2 - y_1| + |z_2 - z_1|$$

	1	2	3	4	5	6	7	8	9	10
1	0	10	3	4	16	32	19	14	15	
2	10	0	11	10	16	8	22	18	6	7
3	3	11	0	15	15	18	18	13	14	
4	4	10	3	6	16	32	19	14	15	
5	4	10	5	6	0	16	28	19	14	15
6	16	8	15	16	16	0	24	7	4	5
7	32	22	33	32	28	24	0	21	26	27
8	19	13	18	19	19	7	21	0	11	12
9	14	6	13	14	14	4	26	11	0	1
10	15	7	14	15	15	5	27	12	1	0

$$\text{dist}[1][2] = (2-1) + (22-23) + (48-48) \\ = 9 + 1 + 0 = 10.$$

Given $k=3$, Now k^{th} neighbours of each point,

1 → 10 ∵ (3, 4, 4, 10, 16, 15, 18, 32)

2 → 8

3 → 11

4 → 6

5 → 6

6 → 7

7 → 24 → Outlier because $24 >$ threshold

8 → 12

9 → 6

10 → 7

Note:—

Another way of detecting outliers:—

- consider top ' k ' highest distances from k^{th} neighbour distances.

∴ For above examples:—

24, 12, 11 are outliers.

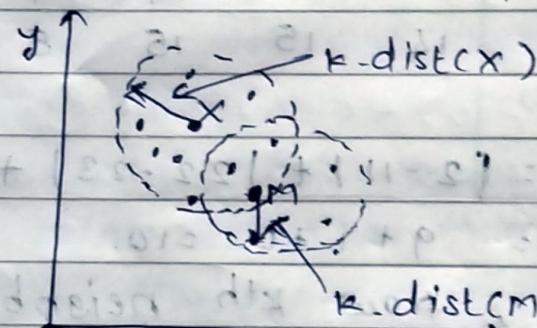
- * Local Outlier Factor (LOF) :- (unsupervised)
 - density based technique of outlier in ml
 - LOF compares the local density of point with the densities of its neighbors.
 - points that have substantially lower density than their neighbors are considered to be outliers.

1) Reachability-distance :-

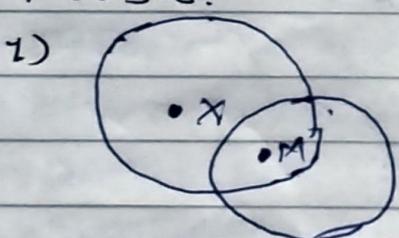
$$\text{reach-dist}_k(M, x) = \max\{\text{k-dist}(M), \text{dist}(M, x)\}$$

Here,

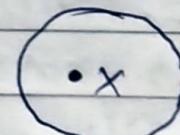
$\text{dist}(M, x)$ = distance of x from M .



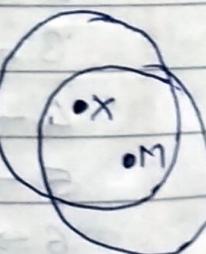
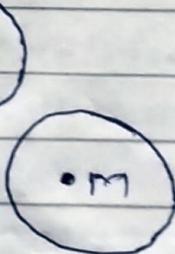
* cases:-



2)



3)



blankness it is the second situation

2) Local Reachability density (LRD):-

$$\text{LRD}_k(M) = \frac{1}{\sum_{x \in N_k(M)} \text{reach-dist}(x)}$$

$$\left[\frac{\sum_{x \in N_k(M)} \text{reach-dist}(x)}{|N_k(M)|} \right]$$

$$3) LOF_{k(M)} = \frac{1}{|N_{k(M)}|} \cdot \sum_{x \in N_{k(M)}} \frac{LRD_K(x)}{LRD_K(M)}$$

* Example :-

index	x	y	
0	2	3	
1	4	5	
2	1	2	
3	2	4	
4	3	4	
5	7	5	
6	3	5	
7	8	9	
8	2	1	
9	3	6	

use distance matrix given in ex. pdf. [k = 2]
finding k-dist & N_k.

index	k-dist	N_k(M)	N_k(M)
0	1.41	(2, 3, 4)	3
1	1.41	(4, 6, 9)	3
2	2.24	(0, 3, 8)	3
3	1.41	(0, 4, 6)	3
4	1.41	(0, 1, 3, 6)	4
5	4.34	(1, 6)	2
6	1.41	(1, 3, 4, 9)	4
7	5.66	(1, 5)	2
8	2	(0, 2)	2
9	1.41	(1, 6)	2

* Reachability distance:-

$$\text{reach_dist}_k(O, 2) = \max(k\text{-dist}(2), \text{dist}(O, 2))$$

$$= \max(2.24, 1.41)$$

$$= 2.24$$

$\therefore LRD_k(O) :-$

x	$k\text{-dist}(x)$	$d(x, M)$	$\text{reach_dist}(M, x)$
2	2.24	1.41	2.24
3	1.41	1	1.41
4	1.41	1.41	1.41

$$\text{Now, } LRD_k(O) = \frac{1}{3}$$

$$\frac{\text{reach_dist}(O, 2) + \text{reach_dist}(O, 3) + \text{reach_dist}(O, 4)}{3}$$

$$= \frac{1}{3}$$

$$= \frac{2.24 + 1.41 + 1.41}{3}$$

$$LRD_k(2) =$$

x	$k\text{-dist}(x)$	$d(x, M)$	$\text{reach_dist}(M, x)$
0	1.41	1.41	1.41
3	1.41	2.24	2.24 2.24
8	2	1.41	2

$$LRD_k(2) =$$

$$\frac{1.41 + \cancel{2} + 2.24 + 2}{3}$$

$$= \cancel{0.62}$$

$$= 0.53$$

LRDK(3) :-

x is k-dist(x) dist. of $d(x, m)$ from each dist (m, x).

→ no. of sq at 1st pric. w/ 1st possibl. 1.41

4. 1.41 (one atm. prob. 1st) 1.41

6. 1.41 1.41 1.41

LRDK(3) ad 1st brin. 1.41 + 1.41

$$\underline{1.41 + 1.41 + 1.41}$$

prob. of 1st sq = 3/27 = 1/9

prob. of 2nd sq = 0.70 (based dist. 1st)

LRDK(4) :-

no. of k-dist(x) dist. of $d(x, m)$ from each dist (m, x)

0. 1st sq. dist. 1.41 (1st possibl. 1.41) 1.41

1. 1st sq. dist. 1.41 (1st possibl. 1.41) 1.41

2. 1st sq. dist. 1.41 (1st possibl. 1.41) 1.41

3. 1st sq. dist. 1.41 (1st possibl. 1.41) 1.41

4. 1st sq. dist. 1.41 (1st possibl. 1.41) 1.41

LRDK(4) = $\frac{1.41 + 1.41 + 1.41 + 1.41}{4}$

4

$$= 0.70$$

$$LOF_k(0) = \frac{1}{|N_k(m)|} \cdot \sum_{x \in N_k(m)} \frac{\text{LRDK}(x)}{\text{LRDK}(m)}$$

$$\frac{1}{3} \times \left(\frac{\text{LRDK}(2)}{\text{LRDK}(0)} + \frac{\text{LRDK}(3)}{\text{LRDK}(0)} + \frac{\text{LRDK}(4)}{\text{LRDK}(0)} \right)$$

$$= \frac{1}{3} \times \left(\frac{0.53}{0.59} + \frac{0.70}{0.59} + \frac{0.70}{0.59} \right)$$

$$= 1.09$$

* Isolation Forest (iForest) :- unsupervised.

- main idea is to isolate anomalies by creating binary trees & using them to partition the dataset into smaller subsets.

* Steps:-

- Branching of tree starts by selecting a random feature first. And then branching is done on random threshold.
- If value is less than threshold, it goes to the left branch else to right branch.

* Clustering :- unsupervised.

- task of grouping data points based on their similarity with each other.

- aims at forming group of homogenous data points from a heterogeneous dataset.
- It evaluates the similarity based on a metric like Euclidean distance, Manhattan distance etc. & then group the points with highest similarity score together.

clustering.

Hierarchical clustering

- 1) Agglomerative
- 2) Divisive

K-Means.

- 1) k-medoids
- 2) other variations

Example :— $k=2$, cluster them with centroids $(2, 1)$ & $(4, 5)$. Determine new centroids.
Use euclidian distance.

euclidian distance = $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$
 $c_1 = \text{cluster } (2, 1)$ $c_2 = \text{cluster } (4, 5)$ (min(c_1, c_2)

x	y	$c_1(2, 1)$	$c_2(4, 5)$	cluster.
2	3	$\sqrt{(2-2)^2 + (3-1)^2} = 2$	2.82	c_1
4	5	4.47	0	c_2
1	2	1.41	4.24	c_1
2	4	3.61	2.28	c_2
8	4	3.16	1.41	c_2
7	5	6.40	3	c_2
3	5	4.12	1	c_2
8	9	10	5.65	c_2
2	1	0	4.47	c_1
3	6	5.09	1.41	c_2

$(2, 1)$

$$\therefore \sqrt{(2-2)^2 + (3-1)^2} = \sqrt{(2)^2} = 2$$

$$\sqrt{(2-4)^2 + (3-5)^2} = \sqrt{4+4} = \sqrt{8} = 2.82$$

$$\therefore c_1 = (2, 3), (1, 2), (2, 1) \quad (2, 8)$$

$$c_2 = (4, 5), (2, 4), (3, 4), (7, 5), (3, 5), (8, 9) \\ (3, 6)$$

Now, new centroids:-

For c_1 :

$$x = \frac{\sum x_i}{N} = \frac{2+1+2}{3} = 1.67$$

$$y = \frac{\sum y_i}{N} = \frac{3+2+1}{3} = 2$$

$$\therefore \text{centroid} = (1.67, 2)$$

For C2:

$$\bar{x} = \frac{\sum x_i}{n} = \frac{2+3+7+3+8+3+4}{7} = 4.29.$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{4+4+5+5+9+6+5}{7} = 5.43.$$

$$\therefore \text{centroid} = (4.29, 5.43)$$

* Agglomerative algorithm:-

Given distance matrix:-

	1	2	3	4	5
1	0				
2	9	0			
3	8	7	0		
4	6	5	9	0	
5	11	10	2	8	0

smallest-distance high-similarity. (2)

combine (3,5).

New matrix,

(3,5)	1	2	4	11
(3,5)	0	(1,2)	(1,4)	(2,4)
1	9	0		
2	10	9	0	
4	9	6	5	0

$\text{dist}(i, (3,5)) := \min_{j \in \{1, 2, 4\}} \text{dist}(i, j)$

$\text{dist}(1, 3) = 3 \quad \text{dist}(1, 5) = 11$

$(3,11)$, we have 3 strategies to choose

1) single linkage = $\min(3, 11) = 3$

2) complete linkage = $\max(3, 11) = 11$

3) Average linkage = $\text{Avg}(3, 11) = 7$

Throughout the problem, we will use complete linkage.

$\text{dist}(2, (3, 5))$:-

$$\text{dist}(2, 3) = 7 \quad \max(7, 10) = 10.$$

$$\text{dist}(2, 5) = 10$$

$\text{dist}(4, (3, 5))$:-

$$\text{dist}(4, 3) = 9 \quad \max(9, 8) = 9.$$

$$\text{dist}(4, 5) = 8$$

Next smallest distance = 5.

combine $(2, 4)$

		<u>(3, 5)</u>	<u>(2, 4)</u>	1	
<u>(3, 5)</u>	0	8	2	4	5
<u>(2, 4)</u>	10	0			
1	11	<u>9</u>	0		

$\text{dist}((2, 4), (3, 5))$

$$\text{dist}(2, 3) = 7 \quad \text{dist}(2, 5) = 9$$

$$\text{dist}(2, 5) = 10 \quad \text{dist}(4, 5) = 8$$

$$\max = 10.$$

$\text{dist}(1, (2, 4))$:

$$\text{dist}(1, 2) = 9 \quad \text{dist}(1, 4) = 6$$

$$\max = 9.$$

New smallest dist = 9.

combine $(1, (2, 4))$

(3, 5) (1, 2, 4)

<u>(3, 5)</u>	0	
<u>(1, 2, 4)</u>	11	0

$\text{dist}((3, 5), (1, 2, 4))$:-

$$\text{dist}(3, 1) = 3$$

$$\text{dist}(5, 1) = 11$$

$$\text{dist}(3, 2) = 7$$

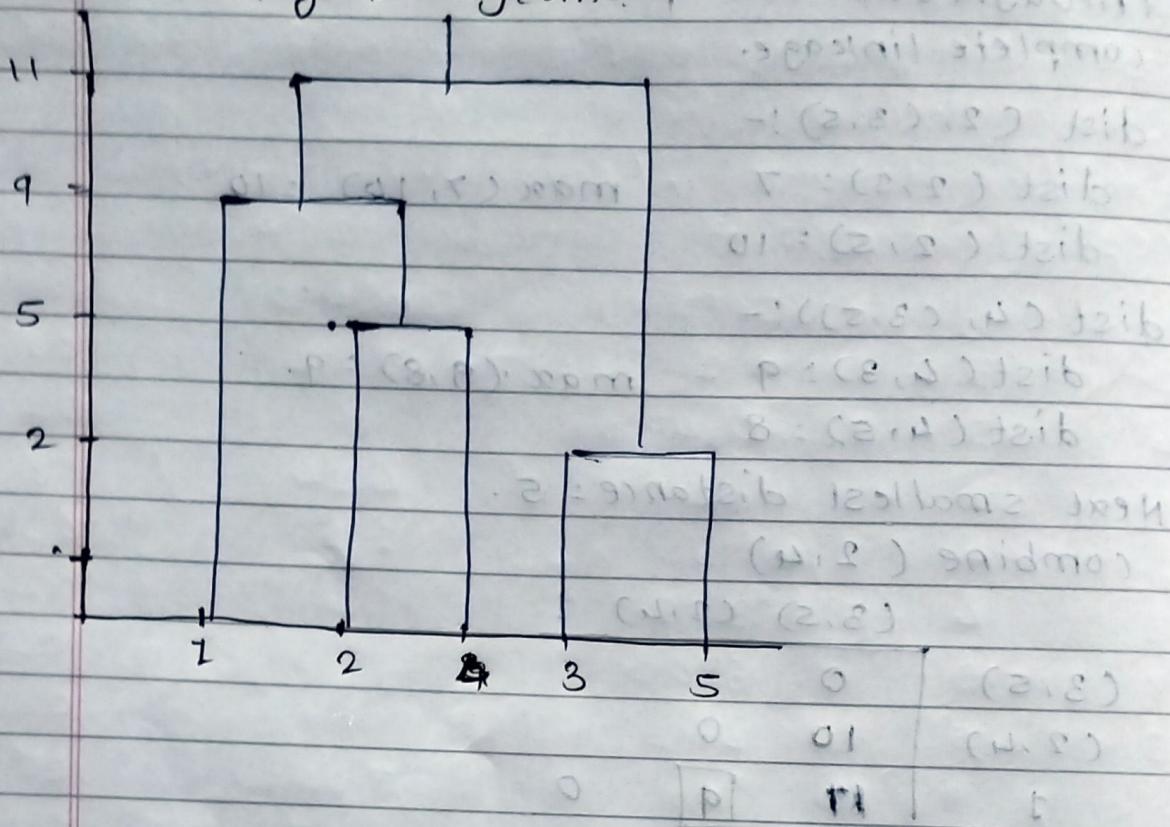
$$\text{dist}(5, 2) = 10$$

$$\text{dist}(3, 4) = 9$$

$$\text{dist}(5, 4) = 8$$

$$\max = 11$$

Drawing Dentogram - with two deciduous



((2, 8), (4, 2)) t2ib
 P = (2, 2) t2ib 5 = (2, 2) t2ib
 8 = (2, 4) t2ib 01 = (2, 8) t2ib
 01 = xpm

:((N, 2), 1) t2ib
 2 = (N, 1) t2ib P = (2, 1) t2ib
 p = xpm

P = t2ib 1201002 wsh
 ((N, 2), 1) snidmos
 (N, 2) (2, 2)

0 (2, 8)
 0 11 (N, 8)

-i((N, 2), 1) (2, 2) t2ib
 11 = (N, 2) t2ib 8 = (2, 2) t2ib
 01 = (2, 2) t2ib 5 = (2, 2) t2ib
 8 = (2, 2) t2ib P = (2, 2) t2ib

1201002