

# Problem A

## Similarity Range Queries

Input File: *testdata.in*  
Time Limit: 10 seconds

### Problem Description

*Similarity range search* has been used in many scientific and industrial applications such as information retrieval, image data analysis, and time-series analysis, etc. Particularly, a similarity range query returns all the data points in a DB that are similar to a query object. Formally, given a data set  $D$  containing data points in  $d$ -dimensional space and a  $d$ -dimensional query object  $q$ .  $q$  retrieves all objects  $o \in D$ , such that  $\text{dist}(q, o) \leq \delta$ , where  $\text{dist}(\cdot, \cdot)$  is the Euclidean distance and  $\delta \in \mathbb{Z}$  a user-specified similarity threshold.

For example, in a image database, every image can be represented by  $d$  features (i.e., dimensions), such as RGB histograms, texture, and shape, etc. Therefore, each image can be viewed as a  $d$ -dimensional data point. In a image retrieval application, a user can specify a  $d$ -dimensional feature vector as a query object and issue the query to the image database to obtain all images within  $\delta$  distance from the query object.

In the following,  $\langle \text{vec}, \delta \rangle$  is used to represent a query object, where  $\text{vec}$  is a  $d$ -dimensional feature vector and  $\delta$  is a user-specified similarity threshold. For example,  $q = \langle [23, 25, 23], 4 \rangle$  is a 3-dimensional query object in which  $\text{vec} = [23, 25, 23]$  and  $\delta = 4$ .

Figure 1 shows 10 3-dimensional data points (i.e.,  $d = 3$ ). After processing  $q = \langle [23, 25, 23], 4 \rangle$ , data points  $p_1$ ,  $p_4$  and  $p_9$  are returned by the database, because the distance of these data points to  $q$  is within  $\delta = 4$ <sup>1</sup>.

---

<sup>1</sup> $\text{dist}(q, p_1) = \sqrt{(23 - 22)^2 + (25 - 24)^2 + (23 - 24)^2} = 1.732$ .

ID	content
$p_1$	[22, 24, 24]
$p_2$	[27, 29, 22]
$p_3$	[20, 30, 21]
$p_4$	[23, 23, 23]
$p_5$	[23, 20, 23]
$p_6$	[25, 20, 22]
$p_7$	[23, 26, 30]
$p_8$	[30, 23, 27]
$p_9$	[24, 24, 22]
$p_{10}$	[26, 20, 23]

Figure 1: The content of a database.

## Technical Specifications

1. The number of data points in a database is smaller than or equal to 2,000,000.
2. The dimensionality  $d = 3$ .
3.  $1 \leq \delta \leq 5$ .
4. The attribute value of each dimension is an integer within the range  $[0, 100]$ .
5. The number of queries in a test case is smaller than or equal to 300.

## Input Format

The input consists of two parts. The first part is the content of the database, and the second part contains several test cases. The two parts are divided by a line containing the integer  $-1$ .

The first part consists of  $N$  data points ( $1 \leq N \leq 2,000,000$ ). Each data point is a single line containing  $d$  integers delimited by a comma. The  $i$ -th ( $i = 1, 2, \dots, d$ ) integer indicates the  $i$ -th attribute value of the data point.

The second part contains several test cases. The first line of the second part is an integer that indicates the number of test cases. Each test case contains several similarity range queries. The two test cases are divided by a line containing the integer  $-1$ . A similarity range query is represented by two lines. The first line is an integer that represents the value of  $\delta$ . The

second line contains  $d$  integers delimited by a comma, which is the feature vector of the similarity range query.

The second part of the input is terminated by a line containing the integer  $-1$ .

## Output Format

We denote the answer of query  $q_i$  by  $ans(q_i)$ . That is,  $ans(q_i)$  is set of data points within  $\delta$  distance of  $q_i$  (i.e.,  $ans(q_i) = \{p_i | dist(p_i, q_i) \leq \delta\}$ ). Therefore,  $|ans(q_i)|$  is the cardinality of  $ans(q_i)$ . For each test case, determine  $ans(q_i)$  for each query  $q_i$  and output the maximum  $|ans(q_i)|$  among all queries. That is, output  $\text{MAX}(|ans(q_1)|, |ans(q_2)|, \dots, |ans(q_n)|)$  in a line. For example, in the following sample input, the first test case contains two similarity range queries:  $q_1 = \langle [23, 25, 23], 4 \rangle$  and  $q_2 = \langle [21, 25, 23], 2 \rangle$ .  $ans(q_1) = \{[23, 23, 23], [22, 24, 24], [24, 24, 22]\}$ .  $ans(q_2) = \{[23, 23, 23], [22, 24, 24]\}$ . Therefore,  $\text{MAX}(|ans(q_1)|, |ans(q_2)|) = \text{MAX}(3, 2) = 3$ .

## Sample Input

```

22,24,24
27,29,22
20,30,21
23,23,23
23,20,23
25,20,22
23,26,30
30,23,27
24,24,22
26,20,23
-1
3
4
23,25,23
2
21,25,23
-1
4
21,29,30

```

3  
30,24,29  
-1  
2  
26,20,26  
4  
20,26,21  
-1

### Sample Output

3  
1  
1