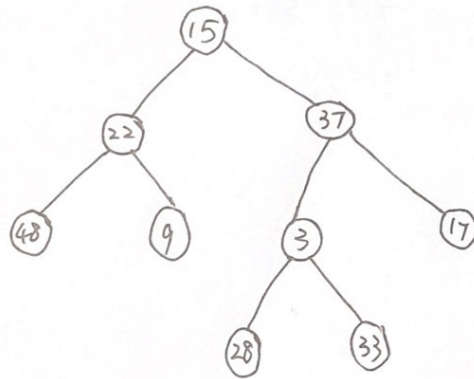Q1. I.



II. We can use binary ~~seach~~ search ~~to~~ examine the midpoint recursively to see if it's larger or smaller than $_\wedge$the element ~~are equal to~~ on the same position in array $[1, 2, ..., n]$. If the element is larger, meaning the missing element is in the left, otherwise it should be in the right half.$_\bullet$ of the array. The pseudo-code is as follows.

```
int findmissing ( int A[], int n ) {
    int b = 0, e = n-1;
    while ( s < e ) {
        int mid = ~~(2)u~~ (b+e) / 2;
        if ( A[mid] < mid + 1) {
            s = mid + 1; }
        else {
            e = mid - 1; }
    }
    if ( A[s] > s+1) { return s+1; }
    ~~else if (A[s] = s+1 return~~
    else { return A[s] + 1; }
}
```

III. Apply the dynamic programming method. ~~For each node.~~

Q2. I. Note that 1 page can hold $\dfrac{2K}{4} = \dfrac{2^{11}}{2^2} = 2^9$ page entries.

The 32-bit address space has $\dfrac{2^{32}}{2^{11}} = 2^{21}$ pages in total.

Thus we need 3 levels of page tables, with one level of page table containing $2^3$ entries, while the other 2 levels of pages tables containing $2^9$ page entries each.

II. 1, 2, 1, 2, 3, 2, 2, 4, 5, 4, 2, 5, 4, 2, 5, 4.

III. (1). abc should be set to 1.

(2) xyz can be set to 0.

(3). The disadvantages ~~is that~~ are that the spin-wait is a waste of CPU time, ~~so~~ and some process may wait forever, which may cause starvation.

To avoid the stravation we can set a queue to hold the processes and ~~plus~~ run them in order.

I. a    $P \leftarrow \Pi_{PID, Category} (\sigma_{category = "Toy"} (PRODUCT))$

    $C \leftarrow \Pi_{CID, Address} (\sigma_{address = "shatin"} (CUSTOMER))$

    $Result \leftarrow \Pi_{EID, Name} (TRANSACTION \bowtie EMPLOYEE \bowtie P \bowtie C)$.

b.    ~~$((EID \; G \; sum (price * Quantity) (TRANSACTION)) \bowtie CUSTOMER)$~~

    $R \leftarrow EID \; G \; sum (price * Quantity)$ as money $(TRANSACTION)$

    $Result \leftarrow \Pi_{EID, NAME, money} (R \bowtie ~~\text{EID}~~ EMPLOYEE)$.

II.  a.  select  EID, EMPLOYEE.Name

    from    TRANSACTION, EMPLOYEE, CUSTOMER, PRODUCT

    where    TRANSACTION.EID = EMPLOYEE.EID

      and  TRANSACTION.PID = PRODUCT.PID

      and  TRANSACTION.CID = CUSTOMER.CID

      and  CUSTOMER.Address = "Shatin"

      and  PRODUCT.Category = "Toy".

  b.  Select  T.EID, E.Name, T.money

    from    (select  EID, sum (Price * Quantity) as money

        from TRANSACTION

        group by EID ) as T, EMPLOYEE as E

    where    T.EID = E.EID.

Q3. III. a. Note that $N_{employee} = 2000$, $B_{employee} = \dfrac{2,000}{100} = 20$.

$N_{transaction} = 100,0000$, $B_{transaction} = \dfrac{100,0000}{50} = 20,000$

For the nested-loop join, if use EMPLOYEE as the outer relation.

we need. ~~Bemployee to~~ $N_{employee} \times B_{transaction} + B_{employee} = 40,000,020$.

block accesses; while using TRANSACTION as the outer

relation. requires $N_{transaction} \times B_{employee} + B_{transaction} = 20,020,000$

block accesses.

Thus we should use TRANSACTION as the outer relation.

For the block nested-loop join:

EMPLOYEE outer: ~~Bemploye~~ $B_{employee} + B_{transaction} = 20,020$.

since ~~the~~ $B_{transaction} > M = 50$.

Thus EMPLOYEE should be chosen as the outer relation.

Q3 III. (b) Nested-loop join:

$$N_{transaction} \times B_{employee} + B_{transaction} = 20,020,000$$

~~Nested-lo~~

Block nested-loop join:

$$B_{employee} + B_{transaction} = 20,020.$$

Q4. I. (a). $\{a\}$ : sup = 5, $\{b\}$ : sup = 4, $\{c\}$ : sup = 4, $\{e\}$ : sup = 3

$\{a,b\}$ : sup = 4, $\{a,c\}$ : sup = 4, $\{a,e\}$ : sup = 3, $\{b,e\}$ : sup = 3

~~$\{a,b,c\}$ : sup = 3, $\{a,b,e\}$ : sup = 3,~~

$\{b,e\}$ : sup = 3,

$\{a,b,c\}$ : sup = 3, $\{a,b,e\}$ : sup = 3.

(b). $\{a\} \rightarrow \{b\}$ : conf $(\{a\} \rightarrow \{b\}) = \dfrac{sup(\{a,b\})}{sup(\{a\})} = 0.8$

$\{a\} \rightarrow \{c\}$ : conf $(\{a\} \rightarrow \{c\}) = 0.8$

$\{b\} \rightarrow \{c\}$ : conf $(\{b\} \rightarrow \{c\}) = 0.75$.

$\{b\} \rightarrow \{e\}$ : conf $(\{b\} \rightarrow \{e\}) = 0.75$

$\{b\} \rightarrow \{a,c\}$ : conf $(\{b\} \rightarrow \{a,c\}) = 0.75$

$\{b\} \rightarrow \{a,e\}$ : conf $(\{b\} \rightarrow \{a,e\}) = 0.75$.

$\{c\} \rightarrow \{a\}$ : conf $(\{c\} \rightarrow \{a\}) = 1.0$.

$\{c\} \rightarrow \{b\}$ : conf $(\{c\} \rightarrow \{b\}) = 0.75$

$\{c\} \rightarrow \{a,b\}$ : conf $(\{c\} \rightarrow \{a,b\}) = 0.75$

$\{e\} \rightarrow \{a\}$ : conf $(\{e\} \rightarrow \{a\}) = 1.0$

$\{e\} \rightarrow \{b\}$ : conf $(\{e\} \rightarrow \{b\}) = 1.0$.

$\{e\} \rightarrow \{a,b\}$ : conf = 1.0.

$\{a,b\} \rightarrow \{c\}$ : conf = 0.75, $\{a,b\} \rightarrow \{e\}$ : conf = 0.75.

$\{a,c\} \rightarrow \{b\}$ : conf = 0.75, $\{a,e\} \rightarrow \{b\}$ : conf = 1.0.

$\{b,c\} \rightarrow \{a\}$ : conf = 1.0, $\{b,e\} \rightarrow \{a\}$ : conf = 1.0.

Q4. (II). Model overfitting is ~~such~~ a case in which the model is ~~too~~ too complicated than the training data, thus the model can very accurately fit to the training data without the ability of generalization.

Method 1 : Improve the size of the training data, or reduce the complexity of the model, (e.g. the ~~number~~ number of ~~par~~ learnable parameters in the model).

Method 2: Introduce regularization term in the loss function ~~to~~ restrict the model to overfit to the training data.

(III). (a) Iter 1 : Allocation : cluster 1 : $\{1\}$

cluster 2 : $\{3, 4, 6, 8\}$

cluster 3 : $\{13, 15, 17, 20, 25\}$

centroids update :   cluster 1 :   $\mu_1 = 1$

cluster 2 :   $\mu_2 = \dfrac{3+4+6+8}{4} = \dfrac{21}{4}$

cluster 3 :   $\mu_3 = \dfrac{13+15+17+20+25}{5} = 18$

Iter 2 :   Allocation :   cluster 1 : $\{1, 3\}$

cluster 2 : $\{4, 6, 8\}$

cluster 3 : $\{13, 15, 17, 20, 25\}$

centroids update:   $\mu_1 = \dfrac{1+3}{2} = 2$.

$\mu_2 = \dfrac{4+6+8}{3} = 6$

$\mu_3 = 18$

Iter 3 :    Allocation =   cluster 1 : $\{1, 3, 4\}$

cluster 2 : $\{6, 8\}$

Cluster 3 : $\{13, 15, 17, 20, 25\}$

centroids :    $\mu_1 = \frac{1+3+4}{3} = \frac{8}{3}$
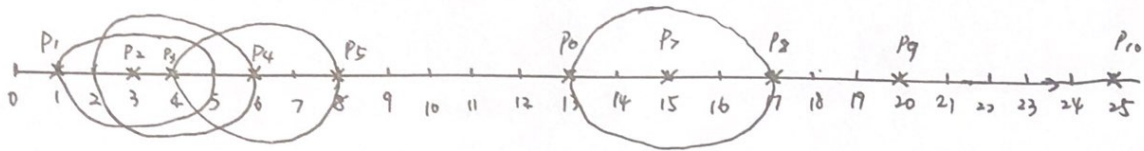
$\mu_2 = \frac{6+8}{2} = 7.$

$\mu_3 = 18$

Iter 4 :    Allocation :   cluster 1 : $\{1, 3, 4\}$

cluster 2 : $\{6, 8\}$

cluster 3 : $\{13, 15, 17, 20, 25\}$

terminates !

III.(b)



core pahts: $P_2$, $P_3$, $P_4$ $P_7$.

border ~~bounder~~ points: $P_1$, $P_5$, $P_6$, $P_8$.

noise points: $P_9$, $P_{10}$.

$P_1$'s  $\varepsilon$-neighbor: $P_2$.
$P_2$'s  $\varepsilon$-neighbor: $P_1$, $P_2$, $P_3$
$P_3$'s  $\varepsilon$-neighbor: $P_2$, $P_3$, $P_4$.
$P_4$'s  $\varepsilon$-neighbor: $P_3$, $P_4$, $P_5$.
$P_6$'s  $\varepsilon$-neighbor: $P_6$, $P_7$, ~~$P_8$~~
$P_7$'s  $\varepsilon$-neighbor: $P_6$, $P_7$, $P_8$
$P_8$'s  $\varepsilon$-neighbor: $P_7$.
$P_9$'s  $\varepsilon$-neighbor: $P_9$
$P_{10}$'s  $\varepsilon$-neighbor: $P_{10}$.

final clusters:   cluster 1: $\{P_1, P_2, P_3, P_4, P_5\}$
                  cluster 2: $\{P_6, P_7, P_8\}$.

Q5.(I)(a) The minimum number of elements in total the method visit is

~~min(3, 4)~~

$$\min(M_1, M_2) + 1.$$

In this case ~~the~~ all ~~the~~ ~~larger~~ elements in the longer sequence are larger than those in the shorter one.

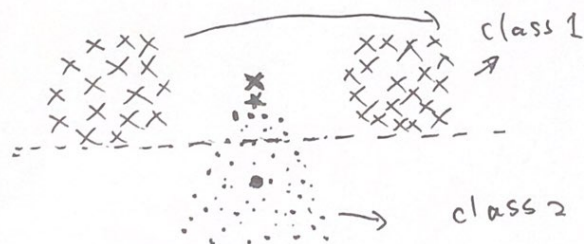Sample : List 1 : $1 \rightarrow 2 \rightarrow 3$

List 2 : $4 \rightarrow 5 \rightarrow 6 \rightarrow 7 \rightarrow 8$.

(b). $M_1 + M_2$, In this case the last element in the shorter list is larger than the last element in the longer one.

Sample : List 1 : $1 \rightarrow 7 \rightarrow 11$

List 2 : $2 \rightarrow 4 \rightarrow 6 \rightarrow 8 \rightarrow 10$.

(II). There are two classes of data below, denoted by cross and dot, respectively. The ~~a~~ bold cross and dot are their corresponding class prototypes. The dotted line is the Rocchio classification boundary. While the bold star "*" is closer to the ~~a~~ centroid of the class 1,



It is ~~more~~ more reasonable to classify it as class 2. In this case the 3-NN method is more suitable.

III. $p(x|c)$ should equal to the fraction of times in which word $x$ appears among all words in documents of class $c$.

$$p(x = w | c = c_j) = \frac{N(x_i = w) + 1}{N(\text{words in } c_j) + N_w}$$

where $N(x_i = w)$ is the number of times $w$ appears in documents of class $c_j$, and $N(\text{words in } c_j)$ is the total number of words in documents of class $c_j$, $N_w$ is the size of the dictionary.