

Advanced Algorithms Problems and Solutions

Mattia Setzu Giorgio Vinciguerra

October 2016

Contents

1	Hogwarts	2
1.1	Solution 1: Preprocessing-then-Dijkstra	2
1.1.1	Pseudo-code	2
1.2	Solution 2: HogwartsDijkstra	3
1.3	Solution 3: BFS-like traversal	4
2	Paletta	5
2.1	Solution 1: Split and count-inversions	5
3	Range updates	7
3.1	Solution 1: Fenwick lazy a-b sums	7
4	Karp-Rabin	11
4.1	Solution 1: Cumulative shift	11
4.1.1	Construction	11
4.1.2	equals(i, j, l)	11
4.1.3	lce(i, j)	12
5	Depth of a node in a random search tree	13
5.1	Recursive balanced proof	13
5.2	Upper bound	14
6	Hashing sets	15
6.1	$k \in S$	15
7	Family of uniform hash functions	16
7.1	Equal probability	16

1 Hogwarts

The Hogwarts School¹ is modeled as a graph $G = (V, E)$ where V is the set of castle's rooms and $E \subseteq V \times V$ is the set of the stairs. Each stair is labelled with the time of appearance and disappearance, and can be walked in both directions, therefore the graph is undirected. The goal is to find, if possible, the minimum amount of time required to go from the first to the last room.

1.1 Solution 1: Preprocessing-then-Dijkstra

Dijkstra is able to find the shortest path in a graph with non-negative weights on its edges. Our main idea is to create a Dijkstra compatible graph through a *normalize* function, then apply Dijkstra to it in order to find the shortest path. The core of the preprocessing is the normalize function which computes traversal times between nodes at a given time *time*:

```
1: function NORMALIZE(from, to, time):
2:    $t = \infty$ 
3:
4:   if  $start[v'] \leq t < end[v']$  then                                ▷ No waiting time
5:      $t = t + 1$ 
6:
7:   else if  $t < start[v']$  then                                       ▷ Waiting time
8:      $t = start[v'] + 1$ 
9:
10:  else
11:     $t = \infty$                                                        ▷ Available time already expired
12:  end if
13: return  $t$ 
14: end function
```

Computational cost: $\Theta(n^2)$ if the vertex set is implemented as an array.
 $O(|E| + |V| \log |V|)$ with Fibonacci heap.

The normalize function is then applied to a node traversal:

1.1.1 Pseudo-code

```
1: create vertex set Q of unvisited nodes
2: create vertexes set E' of edges weight
3: time = 0                                                         ▷ Initial time for traversal
4: edges = stairs_of(0)      ▷ Get the incoming and outgoing edges of the
   starting node
5: function PROCESS(node, time)
```

¹http://didawiki.cli.di.unipi.it/lib/exe/fetch.php/magistraleinformatica/alg2/algo2_16/hogwarts.pdf

```

6:
7:   if  $edge \in visited\_edges$  then return
8:   end if
9:    $traversal\_time = \infty$ 
10:  for each  $neighbor \in neighbors\_of\_node$  do
11:     $traversal\_time = traversal\_time(node, neighbor, time)$ 
12:     $E'[0][node] = traversal\_time$   $\triangleright E'[i][j]$  holds the weight/traversal
13:                                 $\triangleright$  time for the stair between  $i$  and  $j$ 
14:    for each  $new\_neighbor \in neighbors\_of\_neighbor$  do
15:       $normalize(neighbor, new\_neighbor, traversal\_time)$ 
16:    end for
17:  end for
18:   $dijkstra\_graph = \{V, E'\}$ 
19:   $t = dijkstra(dijkstra\_graph)$ 
20:
21:  if  $t == \infty$  then return -1;
22:  elsereturn  $t$ ;
23:  end if
24: end function

```

Computational cost. See the previous section.

1.2 Solution 2: HogwartsDijkstra

```

1: function HOGWARTSDIJKSTRA( $G$ ):
2:   create vertex set  $Q$  of unvisited nodes
3:   for each vertex  $v \in V$  do  $\triangleright$  initialization
4:      $time[v] \leftarrow \infty$   $\triangleright$  unknown time from source to  $v$ 
5:     add  $v$  to  $Q$   $\triangleright$  all nodes initially in  $Q$ 
6:   end for
7:    $time[0] \leftarrow 0$   $\triangleright$  time from source to source
8:   while  $Q \neq \emptyset$  do
9:      $u \leftarrow x \in Q$  with  $\min\{time[x]\}$ 
10:    remove  $u$  from  $Q$ 
11:    for each neighbor  $v$  of  $u$  do:
12:      if  $time[u] \leq appear[v]$  then
13:         $alt \leftarrow appear[v] + 1$   $\triangleright$  wait the appearance of the stair
14:      else if  $time[u] < disappear[v]$  then
15:         $alt \leftarrow time[u] + 1$   $\triangleright$  use the stair
16:      else
17:         $alt \leftarrow \infty$   $\triangleright$  the stair has already disappeared
18:      end if
19:      if  $alt < time[v]$  then
20:         $time[v] \leftarrow alt$   $\triangleright$  a quicker path to  $v$  has been found

```

```

21:         end if
22:     end for
23: end while
24: return time[|N| - 1]
25: end function

```

Computational cost. See the previous section.

1.3 Solution 3: BFS-like traversal

```

1: function REACH( $N, M, A[], B[], appear[], disappear[]$ )
2:   for  $i = 0$  to  $M - 1$  do
3:     edges_ $[A[i]]$ .push_back(make_pair( $i, B[i]$ ))
4:     edges_ $[B[i]]$ .push_back(make_pair( $i, A[i]$ ))
5:   end for
6:   for  $i = 0$  to  $N - 1$  do
7:     done_ $[i] \leftarrow false$ 
8:     distance_ $[i] \leftarrow \infty$ 
9:   end for
10:  reached_ $[0]$ .push_back(0)
11:  distance_ $[0] \leftarrow 0$ 
12:  for  $t = 0$  to  $MAX\_TIME$  do
13:    for  $v \in$  reached_ $[t]$  do
14:      if not done_ $[v]$  then
15:        for edge  $\in$  edges_ $[v]$  do
16:          staircase  $\leftarrow$  edge.first
17:          neighbor  $\leftarrow$  edge.second
18:          time  $\leftarrow$  max(distance_ $[v], appear[staircase]) + 1$ 
19:          if not done_ $[neighbor]$ 
20:             and distance_ $[v] < disappear[staircase]$ 
21:             and time  $<$  distance_ $[neighbor]$  then
22:            distance_ $[neighbor] \leftarrow$  time
23:            reached_ $[time]$ .push_back(neighbor)
24:          end if
25:        end for
26:        done_ $[v] \leftarrow true$ 
27:      end if
28:    end for
29:  end for
30:  return (distance_ $[N - 1] == \infty$ )? - 1 : distance_ $[N - 1]$ 
31: end function

```

Computational cost: $O(m + MAX_TIME)$.

2 Paletta

Paletta ordering² is a peculiar ordering technique: given a 3-tuple of elements, paletta takes the central element as pivot and swaps the two elements right before and next to it. To make an example:

$$\{3, 2, 1\} \xrightarrow{\text{paletta}} \{1, 2, 3\} \quad (1)$$

We now want to develop an algorithm to order any array through paletta ordering with the minimal number of swaps. You should see as not every array can be ordered like this:

$$\{3, , 1\} \xrightarrow{\text{paletta}} \{1, 2, 3\} \quad (2)$$

2.1 Solution 1: Split and count-inversions

We should note that the following properties hold:

1. Every element can be a pivot, but the first and the last one, as they have respectively no elements before and after them.
2. Every element can be swapped as many times as necessary, but only with elements of the same 2-remainder (numbers in even positions can only be swapped with numbers in even positions, the same holds for odd indexes). More formally $|a| = N, \forall i \in [1, N - 2], \neg \exists \text{swap}(i, j) : i \bmod 2 \neq j \bmod 2$
3. The least number of swaps does not backtrack any element. Formally, let k be the minimal number of swaps applied to an array, backtracks included. By hypothesis, k is minimal, but at least $m, m > 0$ backtrack swaps have been operated, therefore we found a $k' = k - m : k' < k$, a new minimal number of swaps: contradiction.

Given item 2, we can split our array in two, even and odd numbers, and order them counting the swaps. In our example we'll use *mergesort*, as it runs in $\log_2 n$, does backtrack elements, and is very well-known. Clearly, given an array, a swap happens when an element is pushed back, pulling the one between its new position and the old one ahead: we can map this behaviour in the merge routine of mergesort: the array merged is able to push back elements from its right pointer to the new array, moving them back of $(m - i) + (j - m)$ positions, where m is the dimension of the current two sub-arrays to merge. Provided that our edited version of mergesort ran successfully on both the even-index and odd-index, we now need to verify if by merging them we obtain an ordered array. Intuitively, the merged array will start with the first element of the even-index arrays, followed by the first of the odd-index array, followed by the second of the even-index array, and so on. To check for these elements is pretty trivial and can be done in linear time. Follows the pseudo-code for the edited version and *snake_check* function:

²http://didawiki.cli.di.unipi.it/lib/exe/fetch.php/magistraleinformatica/alg2/algo2_16/paletta.pdf

```

1: function MERGE_WITH_PALETTE(left, right, k):
2:   ...                                     ▷ merge instructions
3:
4:   if right > left then
5:     paletta_count = paletta_count + 1
6:     ...
7:   end if
8: end function

1: function SNAKE_CHECK:
2:   even, odd = 0
3:   for ;even, odd < N even = even + 1, odd = odd + 1 do
4:
5:     if a[even] > a[odd] then
6:       return -1
7:     end if
8:   end for
9:   return paletta_count
9: end function

```

<p>Computational cost: $\Omega(n \log_2(n)), \Theta(n \log_2(n)), O(n \log_2(n))$.</p>
--

3 Range updates

Consider an array C of n integers, initially all equal to zero. We want to support the following operations:

- **update(i, j, c):** where $0 \leq i \leq j \leq n - 1$ and c is an integer: it changes C such that $C[k] = C[k] + c$ for every $i \leq k \leq j$.
- **query(i)** where $0 \leq i \leq n - 1$: it returns the value of $C[i]$.
- **sum(i, j)** where $0 \leq i \leq j \leq n - 1$: it returns $\sum_{k=i}^j C[k]$.

Design a data structure that uses $O(n)$ space and implements each operation above in $O(\log(n))$ time. Note that $query(i) = sum(i, i)$ but it helps to reason. [Hint to further save space: use an implicit tree such as the Fenwick tree (see wikipedia).]

3.1 Solution 1: Fenwick lazy a-b sums

Let T be a segmented binary tree over a continuous interval $I : [0, N - 1]$ s.t. its leafs are the points in I , and the parent of two nodes comprises of their interval:

$$n' \cup n'' = n, n' \cap n'' = \emptyset \text{ s.t. } n \text{ is the father of } n', n''$$

T will keep track of the prefix sums for every interval. We define a function

$$s' : [0, n - 1] \rightarrow \mathbb{N} \quad (3)$$

that given a node in T returns the value associated with I , namely the cumulative sum of that interval.

In order to reduce the computational cost, we introduce a lazy algorithm that doesn't propagate sums over T as they are streamed in the input, which means $s'(i)$ might not be accurate at a given time t for any of the requested operation.

We'll instead either compute over T or update T as necessary. Let us define a function to do so:

$$l : \mathbb{N} \rightarrow (\mathbb{N} \cup \{\epsilon\}, \mathbb{N}) \quad (4)$$

to keep track of our lazy sums:

$$s(n) = \begin{cases} \epsilon, - & \text{if no lazy prefix sum is in that interval} \\ k, m & \text{if a lazy sum of } k \text{ is to be propagated to } m \end{cases}$$

The *query* function is then trivial:

```

1: function QUERY( $I, i, sum$ ):
2:
3:   if  $I.size == 1$  then                                ▷ Return found value
4:     return  $I.sum$ 
5:   end if
```

```

6:
7:   if lazy(I),  $i \in I.left, i \neg \in I.right$  then                                ▷ Lazy on left child
8:       lazy(I) = False
9:       query( $I.left, i, sum + I.sum$ )
10:   end if
11:
12:   if lazy(I),  $i \in I.right, i \neg \in I.left$  then                                ▷ Lazy on right child
13:       lazy(I) = False
14:       query( $I.right, i, sum + I.sum$ )
15:   end if
16:
17:   if lazy(I),  $i \in I.right, i \in I.left$  then                                ▷ Lazy on both
18:       lazy(I) = False
19:       query( $I.right, i, j, sum + I.sum$ ) + query( $I.left, i, j, sum + I.sum$ )
20:   end if
21:
22:   if !lazy(I),  $i \in I.left$  then                                                ▷ Not lazy on left child
23:       query( $I.left, i, sum$ )
24:   end if
25:
26:   if !lazy(I),  $i \in I.right$  then                                                ▷ Not lazy on right child
27:       query( $I.right, i, sum$ )
28:   end if
29:
30:   if !lazy(I),  $i \in I.right, i \in I.left$  then                                ▷ Not lazy on both
31:       sum( $I.right, i, sum$ )
32:   end if
33: end function
1: function SUM( $I, i, j, sum$ ):
2:
3:   if  $I.size == 1$  then                                                        ▷ Return found value
4:       return  $I.sum + sum$ 
5:   end if
6:
7:   if lazy(I),  $i \in I.left, i \neg \in I.right$  then                                ▷ Lazy on left child
8:       lazy(I) = False
9:       sum( $I.left, i, j, sum + I.sum$ )
10:  end if
11:
12:  if lazy(I),  $i \in I.right, i \neg \in I.left$  then                                ▷ Lazy on right child
13:      lazy(I) = False
14:      sum( $I.right, i, j, sum + I.sum$ )
15:  end if
16:
17:  if lazy(I),  $i \in I.right, i \in I.left$  then                                ▷ Lazy on both

```



```

18:     lazy(I) = False
19:     sum(I.right, i, j, sum + I.sum) + sum(I.left, i, j, sum + I.sum)
20: end if
21:
22: if !lazy(I), i ∈ I.left then                                ▷ Not lazy on left child
23:     sum(I.left, i, sum)
24: end if
25:
26: if !lazy(I), i ∈ I.right then                                ▷ Not lazy on right child
27:     sum(I.right, i, sum)
28: end if
29:
30: if !lazy(I), i ∈ I.right, i ∈ I.left then                    ▷ Not lazy on both
31:     sum(I.right, i, sum)
32: end if
33: end function
1: function UPDATE(I, i, j, k):
2:
3:     if I.size == 1 then                                        ▷ Return found value
4:         return I.val + = update
5:     end if
6:
7:     if lazy(I), i ∈ I.left, i ∉ I.right then                    ▷ Lazy on left child
8:         lazy(I.left) = True
9:         I.left.val = k
10:    end if
11:
12:    if lazy(I), i ∈ I.right, i ∉ I.left then                    ▷ Lazy on right child
13:        lazy(I.right) = True
14:        I.right.val = k
15:    end if
16:
17:    if lazy(I), i ∈ I.right, i ∈ I.left then                    ▷ Lazy on both
18:        lazy(I) = True
19:        I.val = k
20:    end if
21:
22:    if !lazy(I), i ∈ I.left then                                ▷ Not lazy on left child
23:        update(I.left, i, update)
24:    end if
25:
26:    if !lazy(I), i ∈ I.right then                                ▷ Not lazy on right child
27:        update(I.right, i, update)
28:    end if
29:

```

```

30:   if !lazy(I),  $i \in I.right, i \in I.left$  then           ▷ Not lazy on both
31:       update( $I.right, i, update$ )
32:   end if
33: end function

```

4 Karp-Rabin

Given a string $S : |S| = n$, and two positions $0 \leq i < j \leq (n - 1)$, the longest common extension $lceS(i, j)$ is the length of the maximal run of matching characters from those positions, namely: if $S[i] = S[j]$ then $lceS(i, j) = 0$; otherwise, $lceS(i, j) = \max l \geq 1 : S[i \dots i + l - 1] = S[j \dots j + i - 1]$. For example, if $S = \text{abracadabra}$, then $lceS(1, 2) = 0$, $lceS(0, 3) = 1$, and $lceS(0, 7) = 4$. Given S in advance for preprocessing, build a data structure for S based on the Karp-Rabin fingerprinting, in $O(n \ln(n))$ time, so that it supports subsequent online queries of the following two types:

- $lceS(i, j)$: it computes the longest common extension at positions i and j in $O(\log n)$ time.
- $equals(i, j, c)$: it checks if $S[i \dots i + c - 1] = S[j \dots j + c - 1]$ in constant time.

Analyze the cost and the error probability. The space occupied by the data structure can be $O(n \log(n))$ but it is possible to use $O(n)$ space. [Note: in this exercise, a onetime preprocessing is performed, and then many online queries are to be answered on the fly.]

4.1 Solution 1: Cumulative shift

4.1.1 Construction

In order to save computational cycles on checks over ranges we use a similar structure to the one in the range updates: we compute the hashing on the first character in $O(1)$ time, then roll the hash through the $n - 1$ remaining characters through $nO(1)$ operations. We call H this array; we also denote h_k as the function ca^i computing the Rabin-Karp hash of a string s . The reader shall now see that $\exists h^{-1}(s)$: that is, h is invertible in $O(1)$. The entries $h[i] = \sum_{i \in [0, n-1]} (h(i))$ have cumulative hash and the following properties hold:

- $h[s[i]] = (h[i] - h[i - 1]) / a^{-1}, a^{-1} = a^1$
- $h[i..j] - h[k..l] = (h[l] - h[k - 1]) / a^{-1} - (h[j] - h[i - 1]) / a^{-1}, a^{-1} =$ modular inverse

4.1.2 equals(i, j, l)

EQUALS works on cumulative hashes, subtracting them and scaling them accordingly, as our *rabin* function multiplies by an a^i constant.

```

1: function EQUALS( $i, j, length$ ):
2:    $h_i = h[i + length] - h[i - 1]$ 
3:    $h_j = h[j + length] - h[j - 1]$ 
4:    $h^i = h_i / inv(a, i, l)$ 
5:    $h^j = h_j / inv(a, j, l)$ 
   return  $h^i - h^j == 0$ 
```

```

6: end function
7: function INV( $h, k, l$ ): return  $h^{k-l}$ 
8: end function

```

4.1.3 lce(i, j)

LCE works on cumulative hashes, checks the equality on the middle element of the strings and runs recursively on the half with different hashing. We define LCE as an auxiliary function

```

1: function LCE( $i, j, l$ ):
2:   eq = EQUALS( $i, j$ )
3:   if eq then return  $l$ 
4:   else if  $\neg$  EQUALS( $(j-i)/2, (n-j)/2, l$ ) then return EQUALS( $(j-i)/2,$ 
       $(n-j)/2$ )
5:   else  $(j-i)/2+$  return EQUALS( $(j-i)/2, (n-j)/2$ )
6:   end if
7:    $h_i = h[i + length] - h[i - 1]$ 
8:    $h_j = h[j + length] - h[j - 1]$ 
9:    $h^i = h_i / inv(a, i, l)$ 
10:   $h^j = h_j / inv(a, j, l)$ 
      return  $h^i - h^j == 0$ 
11: end function
12: function INV( $h, k, l$ ): return  $h^{k-l}$ 
13: end function

```

5 Depth of a node in a random search tree

A random search tree for a set S can be defined as follows: if S is empty, then the null tree is a random search tree; otherwise, choose uniformly at random a key $k \in S$: the random search tree is obtained by picking k as root, and the random search trees on $L = \{x \in S : x < k\}$ and $R = \{x \in S : x > k\}$ become, respectively, the left and right subtree of the root k . Consider the randomized QuickSort discussed in class and analyzed with indicator variables CLRS 7.3, and observe that the random selection of the pivots follows the above process, thus producing a random search tree of n nodes. Using a variation of the analysis with indicator variables, prove that the expected depth of a node (i.e. the random variable representing the distance of the node from the root) is nearly $2 \log_2(n)$. Prove that the probability that the expected depth of a node exceeds $2 \log_2(n)$ is small for any given constant $c > 1$. [Note: the latter point can be solved after we see Chernoff's bounds.]

5.1 Recursive balanced proof

Let n be the number of nodes in the input list l , $h = \log_2(n)$ the height of a balanced tree over l , $T(p)$ the tree built over the permutation p of pivots, $d(m)$ be the positional distance of a value m of a partition from the median value of the said partition. Then the following holds:

- $height(T) = h \iff |T.left| = |T.right| \pm 1$ Trivially, let r be the root of a 3-nodes partition: then, if the partition is unbalanced, the lesser one will comprise of 0 nodes, while the greater one of 2, which implies that $height(T.right) = 2$.
- $P = \text{pivot}, d(m) = \pm k \implies height(T.left) = height(T.right) \pm k$. Recursively from the previous statement, a partition unbalanced of one element generates subtrees whose levels differ on a factor of 1. By iterating recursively, their subtrees, if unbalanced by 1, will yield one more level difference. Over k unbalanced pivots on a single subtree, at most k levels will be added to h .
- By the previous statement, it follows that $\nexists T, T' : height(T) >= height(T')$, T balanced, T' unbalanced. As stated, let T', T be the unbalanced/balanced tree respectively; let us cheat with T and switch the root pivot with the first element in its subtree. Now, let us prove by contradiction that T can't stay balanced and that its height will increase. By shifting the tree to the left we have deprived $T.right$ of either 0 levels (in case $T.right$ is able to switch every pivot in its tree with its right subtree root, ending with the rightmost leaf in its subtree) or 1, in case no rightmost leaf is present. Therefore $height(T) <= height(T')$.
- The completely unbalanced tree is the tree with the most levels. By taking partitions of size 0 we constantly force, at each level, one subtree to

disappear. Therefore, its level(s) has to be necessarily transferred to its brother. We then have exactly one node per level, therefore n levels.

Behaviour on random permutations Now let us analyse how the tree depth varies according to random pivot selection. We start by applying the 5.1k-distance to a tree T with $n = 3$ nodes. Trivially, $height(T)$ with balanced tree is equal to two. Now, let us pick either the lowest or the greatest pivot possible: the tree is unbalanced towards either the left or the right, but $height(T) = 2$ in both cases. As the reader can see from 5.1, the distance works in absolute value; it is then clear how, at every permutation for a pivot p , out of the n , there are 2 that generate a tree of the same height: $p = d(P) + k, p = d(P) - k$. Given that at every iteration a node x in a completely unbalanced tree T' has a probability of $\frac{1}{n-i}$, we can define the probability of x being a pivot at level l as:

$$P(x_k) = \frac{1}{n-l} \quad (5)$$

Now, in order for x not to be chosen as pivot in the previous $l - 1$ levels we have:

$$P(x_k) = \prod_{k=1}^{l-1} \left(\frac{1}{n-l+1} \right) \quad (6)$$

Given the height of T , the (harmonic) partial series converges to $\ln(n) + 1$. Let us now add a root r s.t. $T'.right = T, T'.left = T$. We now have to consider the mirror case $\ln(n') + \ln(n')$, given by the previous $n' = n/2$ in the logarithm, since the number of nodes doubled, the $+1$ removed for both, since now neither of $T'.left, T.right$ is the root, and a $+1$ added since a new level has been added.

5.2 Upper bound

By hypothesis,

$$E[d(x) > 2c \ln(n)] <<< 1 \quad (7)$$

By definition the ancestor of a node i are independent random variables, and we can apply the Chernoff bounds over the set $x : d(x) \geq 2 \ln(n)$ of random variables determining the expected distance of nodes.

$$P[X \geq cE[X]] < e^{-c \ln(\frac{c}{e})E[X]}$$

Let us consider $X = 1 \forall i == \ln(n)$, the expected depth of $\ln(n)$, then

$$P[X \geq c \ln(n)] < e^{-c \ln(\frac{c}{e}) \ln(n)}$$

6 Hashing sets

Your company has a database $S \subseteq U$ of keys. For this database, it uses a randomly chosen hash function h from a universal family H (as seen in class); it also keeps a bit vector B_S of m entries, initialized to zeroes, which are then set $B_S[h(k)] = 1 \forall k \in S$ (note that collisions may happen). Unfortunately, the database has been lost, thus only B_S and h are known, and the rest is no more accessible. Now, given $k \in U$, how can you establish if k was in S or not? What is the probability of error? (Optional: can you estimate the size $|S|$ of S looking at h and B_S and what is the probability of error?) Later, another database R has been found to be lost: it was using the same hash function h , and the bit vector B_R defined analogously as above. Using h, B_S, B_R , how can you establish if k was in $S \cap R$ (union), $S \cup R$ (intersection), or $S \setminus R$ (difference)? What is the probability of error?

6.1 $k \in S$

Trivially for $B_S[h(k)] = 0$ we can answer FALSE with $P(\text{error}) = 0$. Let us analyse the opposite case, $B_S[h(k)] = 1$. Let $i \in [0, m]$ be some index s.t. $B_S[i] = 1$, and let $cl_S(i)$ be the list of $k \in S : h(k) = i$ for some set $S \in \mathcal{P}(S)$. We can then denote the sets $cl_U := \{k_U : k_U \in U, h(k) = i\}$, $cl_S := \{k_S : k_S \in S, h(k) = i\}$; it is trivial to show that

- $cl_S \subseteq cl_U$ as $S \subseteq U$.
- $|cl_S(k)| \leq |cl_U(k)| \forall k \in U$ as $S \subseteq U$.

Let us not try and estimate $|cl_S(k)|$: given h is universal, we have an expected value of collisions of $E[X_k] \approx \frac{1}{m} \forall k \in S$, that is

$$\begin{aligned} P(h(k^0) = c) &= \frac{1}{m} \\ P(h(k^1) = c) &= \frac{1}{m^2} \\ P(h(k^{i-1}) = c, h(k^i) = c) &= \frac{1}{m^i} \\ P(h(k) = c, \forall k \in S) &= \frac{1}{m^{|S|}} \end{aligned}$$

We can similarly compute the probability of not collision by simply replacing $\frac{1}{m}$ with $(1 - \frac{1}{m})$: $P(h(k) \neq c, \forall k \in S) = 1 - \frac{1}{m^{|S|}}$.

Given our estimate of the collision list, we can now compute an estimate of the error probability: by 6.1, we give an erroneous answer whenever $k \in cl_U(k), k \notin cl_S(k)$, that is we have a margin of error of $cl_U(k) \setminus cl_S(k)$ whose size is $|cl_U(k)| - |cl_S(k)|$. Given the set of k for which $h(k) = i$ the *bad answers* are then

$$1 - P(\text{good answer}) = 1 - (1 - \frac{1}{m})^{|S|} \quad (8)$$

7 Family of uniform hash functions

The notion of pairwise independence says that, for any $x_1 \neq x_2, c_1, c_2 \in \mathbb{Z}_p$, we have that

$$\mathbf{P}(h(x_1 = c_1), h(x_2 = c_2)) = \mathbf{P}(h(x_1 = c_1)) * \mathbf{P}(h(x_2 = c_2)) \quad (9)$$

In other words, the joint probability is the product of the two individual probabilities. Show that the family of hash functions $H = h_{a,b}(x) = ((ax + b) \bmod p) \bmod m$: $a \in \mathbb{Z}_p^*, b \in \mathbb{Z}_p^*$ is *pairwise dependent* where p is a sufficiently large prime number ($m + 1 \leq p \leq 2m$).

7.1 Equal probability

By linear algebra, $ak \bmod p = c \forall k \in \mathbb{N}$. If we were to cap $ak \bmod p = c \forall k \in K, K_N = k_i : k_i < N$

$$\mathbf{P}(h(x_i = c_i)) = \frac{1}{m^2} = \mathbf{P}(h(x_1 = c_1), h(x_2 = c_2)) \quad (10)$$

Given x_i we define as $m_i = (ax_i + b) \bmod p$; since $(ax_i + b)$ is a *linear transformation* m_i is unique. It follows trivially that $\mathbf{P}((ax_i + b) = d) = \frac{1}{p}$. Now, the same goes for x_1, x_2 :

$$(ax_1 + b) \equiv d \quad (11)$$

$$(ax_2 + b) \equiv e \quad (12)$$

By the Chinese remainder theorem the above system has only one solution, therefore $\mathbf{P}((ax_1 + b) = d) = \frac{1}{p}$ and e become independent of d . Having

$$\mathbf{P}(h(x_1 = c_1), h(x_2 = c_2)) = \mathbf{P}(h(x_1 = c_1)) * \mathbf{P}(h(x_2 = c_2)) = \frac{1}{p} * \frac{1}{p} = \frac{1}{p^2} \quad (13)$$

Which, since $m + 1 \leq p \leq 2m$, proves the assumption.