# Datastructures and Algorithms GD

**Big-O-Highscores** 

### HVA

16 september 2017 Sjors Gielen (500765899) & Marc Griepspoor (500710937)

# **Datastructures and Algorithms GD**

#### **Big-O-Highscores**

Before any assignment was started Sjors made a simple test runner to keep track of tests in a consistent manner. The test runner looks like this:



The test runner uses a dictionary to easily add in the excerisze classes along with an integer of how many times the test suite is meant to run the code.

An example of the test suite running Excersize\_1 5000 times 3 times would look like this in the console:

Datastruct_and_alg Average time = Average result =	go_excersizes 0,0	Excersize_1 021754520000000 0,59340	milliseconds, 00	All full	test	5000 tests time:	complet	ed 109,79120000	00000000	milliseconds(	includes	upper f	for loop f	or tests)
Datastruct_and_alg Average time = Average result =	go_excersizes 0,0	Excersize_1 021594599999999 0,43546	milliseconds, 00	All full	test	5000 tests time:	complet	ed 108,83820000	30000000	milliseconds(	includes	upper f	for loop f	or tests)
Datastruct_and_alg Average time = Average result =	go_excersizes 0,0	Excersize_1 023314819999999 0,25306	milliseconds, 00	All full	test	5000 tests time:	complet	ed 117,50340000	0000000	milliseconds(	includes	upper f	For loop f	or tests)
press enter key to =	o close													

#### **Question 1:**

Write a program that generates an array with 100 random numbers (high scores between 0 and 10000) and then checks if the numbers are unique - output is yes or no:

For question one we simply made an array of size 100 and used a for loop to iterate through each.

With the data setup we made a number of methods. First we did the straight forward solution of having a double for loop run through it. But we ended up creating a hashmap buffer to

which we would attempt to add new indexes. Once an index was impossible to be added we would have our result. This way the algorithm is O(f(n)),  $\Omega(f(1))$ .

<pre>33 34 35 35 35 36 4 36 37 37 4 38 39 4 39 4 39 4 4 5 4 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5</pre>	32	ΠĿ	public static class HasDupe
<pre>public static int Hasuuplicate(this lenumerable(int&gt; source, bool counting = false) {     int count = 0;     if (source == null)     {         throw new ArgumentNullException(nameof(source));     }      var checkBuffer = new HashSet<int>();     foreach (int item in source)         foreach (int item in source)         if (lcheckBuffer.Add(item))         {             //Console.WriteLine("could not add" + item);             //Console.WriteLine("could not add" + item);             //Console.WriteLine("could not add" + item);             else</int></pre>	33	L	(
<pre>35 36 37 37 38 39 4 39 4 39 4 4 4 4 4 4 4 4 4 4 4 4 4</pre>	34	부는	public static int HasDuplicate(this lenumerable(int> source, bool counting = taise)
<pre>30</pre>	35		
<pre>37 = if (source == null) 38 39 41 40 41 42 42 43 44 44 44 45 5 44 45 6 45 4 46 47 47 48 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4</pre>	36	<u> </u>	int count = 0;
<pre>38 39 39 39 40 40 42 42 42 44 4 4 4 4 4 4 4 4 4 4 4</pre>	37	부	it (source == null)
<pre>39 40 40 4 41 42 42 43 5 6 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5</pre>	38		
<pre>40 41 42 42 43 44 44 44 45 45 45 46 47 47 47 47 47 48 4 47 47 48 4 4 4 4 4</pre>	39		throw new ArgumentNullException(nameof(source));
<pre>41 42 42 43 44 44 44 45 6 45 6 45 6 46 47 47 47 47 47 48 4 4 4 4 4 4 4 4 4 4 4</pre>	40		}
<pre>42</pre>	41		
<pre>43</pre>	42		<pre>var checkBuffer = new HashSet<int>();</int></pre>
<pre>44 45 45 46 47 47 48 48 49 49 49 49 49 49 49 4 50 50 51 51 51 51 52 6 53 53 54 54 54 55 54 55 55 56 5 5 5 5 5 5 5 5</pre>	43	₽!	foreach (int item in source)
<pre>45 E if (!checkBuffer.Add(item)) 46 47 47 48 48 49 49 49 49 49 4 50 50 50 51 51 51 52 53 54 54 54 54 55 56 5 5 5 5 5 5 5 5 5 5 5</pre>	44		
<pre>46 47 47 48 48 49 49 49 49 49 50 50 50 51 52 6 53 52 6 53 54 54 54 55 5 5 5 5 5 5 5 5 5 5 5 5</pre>	45	무는	if (!checkBuffer.Add(item))
<pre>47 48 48 49 49 49 50 50 50 52 6 53 52 53 54 54 54 55 5 5 5 5 5 5 5 5 5 5 5 5</pre>	46		
48       If (counting)         49       {         50       {         51       }         52       Image: state stat	47	ЦĻ	//Console.WriteLine("could not add" + item);
<pre>49 50 50 51 51 52 6 53 53 54 54 55 55 55 56 5 5 5 5 5 5 5 5 5 5 5</pre>	48	믿는	if (counting)
50       count++;         51       }         52       else         53       {         54       return 1;         55       }         56       }         57       }         58       return count;         60       }         61       }	49		
51     }       52     else       53     {       54     return 1;       55     }       56     }       57     }       58     return count;       60     }       61     }	50		count++;
52     □     else       53     {       54     return 1;       55     }       56     }       57     }       58        59     return count;       60     }       61     }	51		
53       {         54       return 1;         55       }         56       }         57       }         58	52	Р.	else
54     return 1;       55     }       56     }       57     }       58	53		
55     }       56     }       57     }       58        59     return count;       60     }       61     }	54		return 1;
56 } 57 } 58 59 return count; 60 } 61 }	55		}
57 } 58 { 59 return count; 60 } 61 {	56		}
58 59 return count; 60 } 61 }	57		}
59 return count; 60 } 61 }	58		
60 } 61 }	59		return count;
	60		
	61		}

#### **Question 2:**

# Execute the program above 5000 times, and report how many times the array contains at least one duplicate. What is the probability that a high score list contains at least one duplicate value?

Mathematically we can define the odds of a number existing twice in the array by doing the following math:

$$chance = 1 - \frac{PossbileNumbers!}{(possibleNumbers)^{indexes} * (possibleNumbers - indexes)!}$$

Where ! is the factorial operation.

Which in this case would be

$$x = 1 - \frac{10000!}{(10000)^{100} * (10000 - 100)!}$$

Resulting in x to approximately be equal to 0.3914.

This value is accurately represented within our test suite. we even decided to run it for 5 million times to get more significance.

Datastruct_and_algo_excersiz	es.Excersize_1	A11	5000 tests	complete	d					
Average time =	0,010312120000000 milliseconds,	full	test time:		52 <b>,</b> 5550000000000000	milliseconds(includes	upper f	or loo	p for	tests)
Average result =	0,401200									
Datastruct and algo excersiz	es.Excersize 1	A11	5000 tests	complete	d					
Average time =	0,008971500000000 milliseconds,	full	test time:		45,9179000000000000	milliseconds(includes	upper f	or loo	p for	tests)
Average result =	0,400000									
Datastruct and algo excersiz	es.Excersize 1	A11	50000000 tests	complete	d					
Average time =	0,008070610313354 milliseconds,	full	test time:	4128		milliseconds(includes	upper f	or loo	p for	tests)
Average result =	0,390553									
Datastruct and algo excersiz	es.Excersize 1	A11	5000 tests	complete	d					
Average time =	0,008030960000000 milliseconds,	full	test time:		41,067000000000000	milliseconds(includes	upper f	or loo	p for	tests)
Average result =	0,390400									
Datastruct and algo excersiz	es.Excersize 1	A11	5000 tests	complete	d	•				
Average time =	0,007845940000000 milliseconds,	full	test time:		40,105600000000000	milliseconds(includes	upper f	or loo	p for	tests)
Average result =	0,386600									
pross opton key to close										
press enter key to crose										

## **Question 3**

Write a program that initializes an array of length 10000 with all values at 0. Then, add +1 to a random element, and repeat 100 times. Now check if the array contains a value greater than 1 (output is a boolean).

For this we created a new file that we would run through our test suite. The new file is as follows:

10	Ð.	class Excersize_1_Q3 : Excersize
11		
12		private static int entities = 100:
12		public sveride int run()
1.5	1	
14		
15		<pre>int[] data = new int[10000];//must be declared in the run method otherwise we get cross contamination</pre>
16		Random random = new Random();
17	ė.	<pre>for(int i = 0; i<entities; i++)<="" pre=""></entities;></pre>
18		{
19		data[random.Next(data.length)]++:
20		· · · · · · · · · · · · · · · · · · ·
21		, J sture Hasters Hasturgeton (data).
21		recuir nashore nashore (uaca),
22		
23		
24		
25	Θ	public static class HasMore
26		{
27	古	public static int HasTwoOrMore(this IEnumerable <int> source, bool counting = false)</int>
28	Ī	{
20		i introumt = 0.
29		in court - 9,
50	무는	if (source == null)
31		
32		throw new ArgumentNullException(nameof(source));
33		
34		var orderedBuffer = source.OrderByDescending( $i \Rightarrow i$ );
35		
36	古	foreach (int item in orderedBuffer)
37	ΠŤ	{
20		if(itam > 1)
20	1	
59		
40	부는	11 (counting)
41		
42		count++;
43		}
44	6	
45		f f
46		return 1:
47		
10	∎ F≦ -	
40	L:	f
49	무는	else if (:counting)//when we are not counting the first index tells is weither or not there are auplicates
50		
51		return 0;
52		}
53		}
54		return count;
55		
56		}
50	• L ·	

As the orderby from LINQ uses stable quicksort the big-O we are dealing with here is  $O(f(n \log(n)))$ .

#### **Question 4**

Execute the program above 5000 times, and report how often at least one duplicate value is found. Again, what is the probability that a high score list contains at least one duplicate value?

	F 1 4 03		5000					
Datastruct_and_algo_excersi: Average time = Average result =	<pre>zes.txcersize_1_Q3 0,985666779999996 milliseconds,</pre>	All full	5000 tests co test time:	ompleted 4931,1431000000000000	milliseconds(includes	upper for	loop for	tests)
Datastruct_and_algo_excersi: Average time = Average result =	zes.Excersize_1_Q3 0,953149839999999 milliseconds, 0,411200	All full	5000 tests co test time:	ompleted 4767,802700000000000	milliseconds(includes	upper for	loop for	tests)
Datastruct_and_algo_excersi: Average time = Average result =	zes.Excersize_1_Q3 0,953563120000001 milliseconds, 0,398600	All full	5000 tests co test time:	ompleted 4769,846400000000000	milliseconds(includes	upper for	loop for	tests)
Datastruct_and_algo_excersi: Average time = Average result =	zes.Excersize_1_Q3 0,952668999999994 milliseconds, 0,429800	All full	5000 tests co test time:	ompleted 4765,407500000000000	milliseconds(includes	upper for	loop for	tests)
Datastruct_and_algo_excersi: Average time = Average result =	zes.Excersize_1_Q3 0,958087619999997 milliseconds, 0,365600	All full	5000 tests co test time:	ompleted 4792,464600000000000	milliseconds(includes	upper for	loop for	tests)
press enter key to close								

The probabilities are found on the average result section of the reports.

#### **Question 5**

Are the probabilities found under 2) and 4) equal? Explain why (not).

In our testing suite we do find the same probability, or at least a close enough estimate over 5 tests.

	cirears reportoatastruct and algo excersizes (batastru	ict and	algo excersizes (biri (bei	bug\Datastruct and algo excersizes.exe
Datastruct_and_algo_excer:	sizes.Excersize_1_Q3	A11	5000 tests	completed
Average time =	0,984631839999999 milliseconds,	full	test time:	4925,770700000000000 milliseconds(includes upper for loop for tests
verage result =	0,352000			
atastruct and algo excer	sizes.Excersize 1 03	A11	5000 tests	completed
Average time =	0,95526980000003 milliseconds,	full	test time:	4778,381200000000000 milliseconds(includes upper for loop for tests
verage result =	0,438400			
atastruct and algo excer	sizes.Excersize 1 03	A11	5000 tests	completed
verage time =	0,952812300000001 milliseconds,	full	test time:	4766,029400000000000 milliseconds(includes upper for loop for tests
verage result =	0,373800			
Datastruct and algo excer	sizes.Excersize 1 Q3	A11	5000 tests	completed
verage time =	0,955386399999999 milliseconds,	full	test time:	4779,04810000000000 milliseconds(includes upper for loop for tests
verage result =	0,379200			
atastruct and algo excer	sizes.Excersize 1 Q3	A11	5000 tests	completed
verage time =	0,965166320000003 milliseconds,	full	test time:	4827,97470000000000 milliseconds(includes upper for loop for tests
verage result =	0,408800			
atastruct_and_algo_excer	sizes.Excersize_1	A11	5000 tests	completed
verage time =	0,007905780000000 milliseconds,	full	test time:	40,405500000000000 milliseconds(includes upper for loop for tests
werage result =	0,390600			
atastruct_and_algo_excer	sizes.Excersize_1	A11	5000 tests	completed
verage time =	0,007794560000000 milliseconds,	full	test time:	39,854900000000000 milliseconds(includes upper for loop for tests
werage result =	0,384600			
atastruct_and_algo_excer	sizes.Excersize_1	A11	5000 tests	completed
verage time =	0,007778400000000 milliseconds,	full	test time:	39,768300000000000 milliseconds(includes upper for loop for tests
werage result =	0,388000			
atastruct_and_algo_excer	sizes.Excersize_1	A11	5000 tests	completed
verage time =	0,00903100000000 milliseconds,	full	test time:	46,217400000000000 milliseconds(includes upper for loop for tests
werage result =	0,397400			
atastruct_and_algo_excer	sizes.Excersize_1	A11	5000 tests	completed
Verage time =	0,008425140000000 milliseconds,	full	test time:	43,099100000000000 milliseconds(includes upper for loop for tests
verage result =	0,385600			
mass onten key to close				

Another interesting note is the huge speed difference here. The original formula is about 100 times faster. This may be related to how the data is setup. Perhaps the test suite should have data be defined outside of the scope of the stopwatch?

More likely thought is our difference in big o notation. The sorting process simply takes up too much time currently.

#### Question 6 & 7

## Consider the program you created to answer question 1 and 3. What is the time complexity (in terms of Big-O) of this program? Motivate your answer.

Time complexities were already handled on questions 1 and 3. For summary they were. For the first solution O(f(n)).

This is the result of simply adding items directly to a hasmap which will be impossible if the hasmap already contains the value, via this we detect a duplicate value. This wat we would only have to iterate once through the entire dataset.

For the second solution  $O(f(n \log(n)))$ .

Due to the sorting process we require the slower O(f(n log(n))). Quicksort may be fast but it doesn't shine when it only gets a small amount of indexes. It reverts back to Selection sort far too fast(instead of binary sort). Keep in mind that n is assumed to be the arraysize here and not the playercount. Increasing the playercount should have no effect on the speed on this algorithm

## Question 8 & 9

#### Taking n as the number of players, fill in the following table for your solution of 1 & 3. Make sure you take the average time over at least 1000 execution runs. Do these results confirm your analysis reported in question 6)?

To make testing this faster and more flexible we modified the scripts to make the variables in question be present in the constructors of the class. This would let us test several configurations in a swift manner. The testing suite received some small modifications to support this as well.

Datastruct_and_algo_excersizes.Excersize_1_03 All 1000 tests completed
Average time = 1,048942900000000 milliseconds, full test time: 1049,84170000000000 milliseconds(includes upper for loop for tests)
Average result = 0,514000 variable n was: 100
Datastruct_and_algo_excersizes.Excersize_1_Q3 All 1000 tests completed
Average time = 0,987760000000000 milliseconds, full test time: 988,20970000000000 milliseconds(includes upper for loop for tests)
Average result = 0,867000 variable n was: 200
Datastruct_and_algo_excersizes.Excersize_1_Q3 All 1000 tests completed
Average time = 1,0421018000000000 milliseconds, full test time: 1042,57130000000000 milliseconds(includes upper for loop for tests)
Average result = 1,000000 variable n was: 400
$Jatas truct_and algo_excersizes.cxcersize_i_U_{0}$ and $1000$ lests completed $1000$ and $1000$ lests completed $1000$ and $1000$ lests completed $1000$ and $1000$
Average time - 1,00500410000000 ministerious, full test time. 1006,50,00000000000 ministerious(includes upper for for tests)
Datastruct and algo excersizes.Excersize 1 All 1000 tests completed
Average time = 0.008382900000000 milliseconds, full test time: 8.57520000000000 milliseconds(includes upper for loop for tests)
Average result = 0,389000 variable n was: 100
Datastruct_and_algo_excersizes.Excersize_1 All 1000 tests completed
Average time = 0,010676100000000 milliseconds, full test time: 10,87940000000000 milliseconds(includes upper for loop for tests)
Average result = 0,855000 variable n was: 200
Datastruct_and_algo_excersizes.Excersize_1 All 1000 tests completed
Average time = 0,013450600000000 milliseconds, full test time: 13,63950000000000 milliseconds(includes upper for loop for tests)
Average result = 1,000000 variable n was: 400
Detectionst and also averaging Exercises 1 All 1000 tests completed
patastruct_min_argo_excersizes.txtersize_1 All 1000 tests Completed
Average time - 0,010/000000000 ministerious, full test time. 10,01/0000000000 ministerious(includes upper for for tests)
aress enter key to close

As expected the Q3 solution does not change in terms of time as the respective n value for it was not altered. More players just means that the chance it returns true increases.

For the Q1 solution though we can see that the time slowly increases. Again this is along with what we would expect. A larger list means the chance the first duplicate is found increases slower.

#### Amendment

After writing the main points through this document we realized that it would currently be far more efficient to remove the sorting from the Q3 solution. In doing so we reduce the big O down to n! - (n-size)! instead of nlog(n). As the probability to find a 2 on the first index goes up with more players.

This means the code has moved to this:



And the eight piece test now looks like:

Also please note that on some of the earlier tests we had the wrong random seed setting in that random parameter of Q3. This results in Q3 having too high probabilities sometimes as the seeds isn't properly being modulated.

This is achieved by doing: Random random = new Random(Guid.NewGuid().GetHashCode()); instead of Random random = new Random();

Datastruct and algo excersi:	zes.Excersize 1 Q3	A11	5000 tests	completed			
Average time =	0,047695080000000 milliseconds,	full t	est time:	239,95990000000000	) milliseconds(includes	upper for loop	for tests)
Average result =	0,382000 variable n	was:	100				
Datastruct_and_algo_excersi	zes.Excersize_1_Q3	A11	5000 tests	completed			
Average time =	0,029778060000000 milliseconds,	full t	test time:	150,18040000000000	<pre>milliseconds(includes</pre>	upper for loop	for tests)
Average result =	0,867800 variable n	was:	200				
Datastruct_and_algo_excersi	zes.Excersize_1_Q3	A11	5000 tests	completed			
Average time =	0,016636180000000 milliseconds,	full t	test time:	84,320000000000000	) milliseconds(includes	upper for loop	for tests)
Average result =	0,999800 variable n	was:	400				
Datastruct_and_algo_excersi	zes.Excersize_1_Q3	A11	5000 tests	completed			
Average time =	0,018489960000000 milliseconds,	full t	test time:	93,54640000000000	<pre>milliseconds(includes</pre>	upper for loop	for tests)
Average result =	1,000000 variable n	was:	800				
Datastruct_and_algo_excersi	zes.Excersize_1	A11	5000 tests	completed			
Average time =	0,009049200000000 milliseconds,	full 1	test time:	46,23860000000000	<pre>milliseconds(includes</pre>	upper for loop	for tests)
Average result =	0,385000 variable n	was:	100				
Datastruct_and_algo_excersi	zes.Excersize_1	A11	5000 tests	completed			
Average time =	0,010438780000000 milliseconds,	full t	test time:	53,15210000000000	<pre>milliseconds(includes</pre>	upper for loop	for tests)
Average result =	0,861000 variable n	was:	200				
Datastruct_and_algo_excersi	zes.Excersize_1	A11	5000 tests	completed			
Average time =	0,013333620000000 milliseconds,	full t	test time:	67,59640000000000	<pre>milliseconds(includes</pre>	upper for loop	for tests)
Average result =	0,999400 variable n	was:	400				
Datastruct_and_algo_excersi	zes.Excersize_1	A11	5000 tests	completed			
Average time =	0,019003700000000 milliseconds,	full t	est time:	95,97870000000000	<pre>milliseconds(includes</pre>	upper for loop	for tests)
Average result =	1,000000 variable n	was:	800				
press enter key to close							

The differences in time is now more in-line with what we expect from our big O notation. The main reason we see Q3 400 be faster than Q3 800 is that the data-setup takes much longer compared on Q3 compared to Q4. A thought is to make the dataset production be handled outside of the stopwatch timers to remove this part of the process from being timed.

Keep in mind that the data creation on both ends is big O = n, therefor we did not take them into account when comparing the actually methods which execute the mutation on the dataset.