

APPLICATION ANALYSIS

HASH FUNCTIONS

PROBLEM DESCRIPTION

Hash tables use a hash function to compute an index into an array of buckets or slots, from which the desired value can be found. It is possible that different elements, say different strings consist of exactly the same hash value. This phenomenon is called collision.

The problem with hash functions is when at least elements that were hashed from a particular data structure, say array hash or rather reference the same value in the hash table. This could cause conflict as well whenever one tries to index the element, in order to return its hash value. A hash function is created which does all the work. There are different types of hash functions. Some hash functions just hash values regardless of whether there might be some duplicates, in terms of hash values. Some hash functions solve collisions, ensuring that each element is hashed to a different place. The limitation of hash tables is that their size is fixed, unlike other related data structures like array lists.

APPLICATION DESIGN AND IMPLEMENTATION

Hash function 1

```
public static int hash1(String key){
    int i =0;
    if(key.equals("")){
        return i;
    }
    else{
        i = 1;
        return i;
    }
}
```

- Hash1 method is of type integer, it takes in a string value, from a certain array. Checks if the string is an empty string and returns a 0 if true, if not it returns 1. Basically this hash function will hash all string values to 1 if they are not empty.

Hash function 2

```
public static int hash2(String key){
    int hashVal = 0;
    for(int i=0;i<key.length();i++){
        hashVal += key.charAt(i);}
    hashVal %= tablesize;
    if(hashVal < 0){
        hashVal+=tablesize;
    }
    return hashVal;
}
```

- This hash function returns an integer value. It takes in a string from a certain array and returns its hashval. It progresses in this by creating a variable hashval which will store the hash value of that string. The function adds the Unicode of all characters present in the string, then takes that Unicode values and apply a modulus of fixed table size, 20011. If the hashval is negative, it is incremented by table size.

Hash function 3

```
public static int hash3(String key){
    int hashVal = 0;
    for(int i=0;i<key.length();i++){
        hashVal = (37*hashVal) + key.charAt(i);
    }
    hashVal%=tablesize;
    if(hashVal<0){
        hashVal+=tablesize;
    }
    return hashVal;
}
```

- Hash3 is a hash function of type integer. The hash function basically calculates hash values by shifting the hash table by a factor of 37, creating some holes or rather empty spaces between each entered element. This is a sufficient hash function in comparison to the first hash function and the second hash fuction. Although collisions may occur, but it reduces the number of collisions in the hash table. The hashVal is multiplied by a factor of 37 and then add Unicode values. A modulus is applied to the hashVal by fixed table size of 20011. The hashVal is returned.

Hash function 4

```
public static int hash4(String key){  
    int prime = 31;  
    int hashVal = 1;  
    hashVal = prime * hashVal + ((key == null) ? 0 : key.hashCode());  
    return Math.abs(hashVal);  
}
```

- hash4 is a hash function of type integer. The hash function takes in string as a parameter, multiply hashvalue by 31, and extract hashcode using hashCode(). It then returns the absolute value of the hashvaue.

Experimental Design and Results

Below is the main method that calculates the expected value or rather the entropy of hashvalues for all hash functions.

```
public static void main(String[] args) throws IOException {
    Hashtable<String,Integer> table = new Hashtable<String,Integer>();
    ArrayList<String> mlist = new ArrayList<String>();
    ArrayList<Integer> mlist2 = new ArrayList<Integer>();
    ArrayList<Integer> hash_values1 = new ArrayList<Integer>();
    ArrayList<Integer> entropy_values2 = new ArrayList<Integer>();
    int index = 0;
    Scanner file = new Scanner(new FileInputStream("testdata"));
    while(file.hasNextLine()){
        String line = file.nextLine();
        String fullname = line.substring((line.lastIndexOf("|")+1);
        table.put(fullname,index);
        mlist.add(fullname);
        index++;
    }

    for(int i =0;i<mlist.size();i++){
        hash_values1.add(hash1(mlist.get(i)));
        //entropy_values1.add(hash1(mlist.get(i)));
        //entropy_values2.add(hash1(mlist.get(i)));
        //System.out.println(entropy_values.get(i));
    }

    //insert number of occurrences of each duplicate
    Map<Integer,Integer> count = new HashMap<Integer,Integer>();
    for (int ind : hash_values1) {
        if (count.containsKey(ind)) {
            count.put(ind, count.get(ind) + 1);
        } else {
            count.put(ind, 1);
        }
    }

    for (Map.Entry<Integer, Integer> entry : count.entrySet()) {
        mlist2.add(entry.getValue()); }

    //entropy calculation
    double partial_entropy = 0;
    for(int j=0;j<mlist2.size();j++){
        double value = (double)mlist2.get(j);
        double probability = value/10000;
        BigDecimal dec_val = BigDecimal.valueOf(probability);
        double formula = -probability*(Math.log(probability));
        BigDecimal final_value2 = BigDecimal.valueOf(formula);
        partial_entropy =Math.abs(partial_entropy+formula);
    }
    BigDecimal final_entropy = BigDecimal.valueOf(partial_entropy);
    if(partial_entropy==0){
        System.out.println("Expected value: "+final_entropy.toPlainString()+" => Certainty");
    }
}
```

```

else{
    DecimalFormat formater = new DecimalFormat("0.00000");

    if(partial_entropy==0){
        System.out.println("Expected value: "+formater.format(final_entropy)+" => Csertainty");
    }

    Else{ System.out.println("Expected value: "+formater.format(final_entropy)+" => Uncertainty");}
    }
    }
    }
}

```

- This is the calculation of the entropy from 10 000 entries of data, specifically 10 000 hashvalues stored in the array list.

Calculation:

```

//entropy calculation
double partial_entropy = 0;
for(int j=0;j<mlist2.size();j++){
    double value = (double)mlist2.get(j);
    double probability = value/10000;
    BigDecimal dec_val = BigDecimal.valueOf(probability);
    double formula = -probability*(Math.log(probability));
    BigDecimal final_value2 = BigDecimal.valueOf(formula);
    partial_entropy =Math.abs(partial_entropy+formula);
}
BigDecimal final_entropy = BigDecimal.valueOf(partial_entropy);
if(partial_entropy==0){
    System.out.println("Expected value: "+final_entropy.toPlainString()+" => Certainty");
}

```

Output:

Hash function 1: entropy = 0.0 => Certainty

Hash function 2: entropy = 6.63181 => Uncertainty

Hash function 3: entropy = 8.89434 => Uncertainty

Hash function 4: entropy = 9.20618 => Uncertainty

Tables

Hash function 1

Hash function 2

FileHomeInsertPage LayoutFormulasDataReviewViewToolsTell me what you want to do

CutCopyPasteFormat PainterClipboard

Calibri11AaWrap Text

BItUFontAlign Center

GeneralConditional FormattingTableStyles

NormalBadGoodNeutral

InsertDelete FormatClear

Sort & Find & Filter & Select

AutoSumFillCellsEditing

D4=COUNTIF(A54:A510003,B4)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	
1																			
2																			
3	Hash Values	Criteria		Mode		Hash values(Sorted)	Mode	Entropy											
4		4655	4655	3			20002	1	8.89434 [from the program]										
5		18255	18255	1			19999	2											
6		2601	2601	2			19999	2											
7		7965	7965	2			19998	1											
8		5087	5087	1			19995	1											
9		15577	15577	2			19994	1											
10		4143	4143	2			19990	1											
11		4306	4306	2			19988	1											
12		31	31	1			19984	1											
13		9713	9713	1			19983	1											
14		17989	17989	1			19982	1											
15		8216	8216	2			19981	1											
16		4944	4944	1			19978	1											
17		16185	16185	1			19977	1											
18		9629	9629	1			19973	1											
19		10383	10383	2			19971	2											
20		6920	6920	1			19971	2											
21		1275	1275	1			19970	1											
22		10964	10964	1			19969	1											
23		17940	17940	1			19967	1											
24		11277	11277	2			19962	1											
25		18201	18201	2			19961	1											
26		3675	3675	3			19958	3											
27		855	855	1			19958	3											
28		16636	16636	2			19958	3											
29		18083	18083	1			19957	1											
30		419	419	1			19955	1											

Sheet1Sheet2Sheet3

8:20 AM 9/15/2017

Hash function 3

Hash function 4

FileEditViewInsertFormatToolsHelp

Calibre

E5

A

B

C

D

E

F

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

Hash Values and entropy(Hash function4)

Hash Values

Criteria

Hash Values(Sorted)

Mode

Entropy

1871882874

1871882874

2146803715

1

9.20616 [from the program]

533997698

533997698

2146371956

1

392791854

392791854

2145358050

1

1718424274

1718424274

2145051074

1

1268398670

1268398670

2144791313

1

1627920968

1627920968

2143633931

1

753001766

753001766

2143476485

1

1787670714

1787670714

2143454810

1

1915585262

1915585262

2143194549

1

970653659

970653659

2143156477

1

1645570865

1645570865

2143132035

1

189532357

189532357

2142973160

1

851895072

851895072

2142945380

1

802950566

802950566

2142645260

1

833796821

833796821

2142642360

1

710144408

710144408

2142619310

1

795646788

795646788

2142547035

1

171400858

171400858

2141790443

1

635347912

635347912

2141237982

1

1775680206

1775680206

2141159348

1

1644957980

1644957980

2140860219

1

1936241121

1936241121

2140854301

1

2004978821

2004978821

2140560006

1

2136611764

2136611764

2140426496

1

336714287

336714287

2140052211

1

Sheet1Sheet4Sheet2Sheet3

PageStyle Sheet4

Average: Sum: 0

Discussion and Conclusion:

The tables above indicate the mode or number of occurrences a hash value of different strings appears. The tables are sorted out in order from smallest to largest. The entropy value was calculated using from a java program. The number of occurrences is thus what we call collision. As the tables show, there are many values that hash to the same place, even though the strings might differ. This is a huge problem because whenever a certain value is indexed from the hashtable it will point to different elements of String and not a unique string.

The entropy as defined before is the expected value of the probability of occurrences. Output for hash function 1 is 0, this is certainty of the outcome while other functions give entropy greater than 1, this represents uncertainty.

In concluding, when the entropy is 0, this means that we are certain about the expected value. From table 1, the hash table consists of only 1s, from the initial index up to 10 000, hence the data won't be random. In general, the entropy is the indication of the randomness collected by the operating system.