

**ComS 535x:**  
**Project Report #2**

Due on Mar 8, 2015

*Instructor: Professor Aduri, Pavankumar*

**Yuanyuan Tang, Chenguang He**

## Contents

<b>1</b>	<b>MinHash</b>	<b>3</b>
1	Your procedure to collect all terms of the documents and the data structure used for this. . . . .	3
2	Your procedure to assign an integer to each term . . . . .	3
3	The permutations used, and the process used to generate random permutations . . . . .	3
4	Brief pseudo code for methods: exactJaccard, minHashSig, approximateJaccard, minHashMatrix . . .	3
<b>2</b>	<b>MinHashAccuracy</b>	<b>6</b>

## MinHash

### 1 Your procedure to collect all terms of the documents and the data structure used for this.

I use Scanner to scan each file in the given folder, and save all modified terms in  $HashSet < String > terms$ . Also, I use  $HashMap < String, HashSet < String >> allStringsPerFile$  to store the file name and its terms. The reason I pick HashSet to store the term is that, when we construct the binary representation of terms, we can have  $O(1)$  access to the terms.

### 2 Your procedure to assign an integer to each term

I make an iterator, iterate all the terms one by one. I start a counter  $j$  at the beginning of iteration. For each term, if it belongs to the file, I put the counter in output array. Therefore, we get an array of a sequence with all terms appeared in increasing order. In here, we use HashSet instead of TreeSet is kind of tricky. Since the iterator of HashSet is not in order, it could result in the different iterator returns after we add new element in HashSet. However, after we read files, the HashSet is final. It will not change, that is why I pick up HashSet instead of TreeSet. (Notice that, TreeSet's iterator is in order and stable).

### 3 The permutations used, and the process used to generate random permutations

I pick a very large prime  $p$  at the beginning,  $p \in [N, 2N]$ , where  $N$  is the number of terms. Then I use  $(a+bx)$  to generate  $p$ , I use prime test from  $N$ . and it randomly comes a number in  $[N, 2N]$ . By the prime number theorem, the average number of prime in  $[N, 2N]$  are  $\frac{N}{\ln(N)}$ , therefore there must exist enough prime numbers in the range.

### 4 Brief pseudo code for methods: exactJaccard, minHashSig, approximateJaccard, min-HashMatrix

**Algorithm 1** Exact Jaccard Calculation

---

```

1: procedure EXACTJACCARD(file1, file2)
2:    $b1 \leftarrow \text{termsBinaryRepresentation of file1}$ 
3:    $b2 \leftarrow \text{termsBinaryRepresentation of file2}$ 
4:   intersection  $\leftarrow 0$ 
5:   union  $\leftarrow 0$ 
6:    $i \leftarrow 0$ 
7:    $j \leftarrow 0$ 
8:   while  $b1[i] \neq -1$  and  $b2[j] \neq -1$  do
9:     if  $b1[i] = b2[j]$  then
10:      intersection  $\leftarrow$  intersection + 1
11:       $i \leftarrow i + 1$ 
12:       $j \leftarrow j + 1$ 
13:     end if
14:     if  $b1[i] < b2[j]$  then
15:        $i \leftarrow i + 1$ 
16:     end if
17:     if  $b1[i] > b2[j]$  then
18:        $j \leftarrow j + 1$ 
19:     end if
20:     union  $\leftarrow$  union + 1
21:   end while
22:   union  $\leftarrow$  union + 1
23:   return  $\frac{\text{intersection}}{\text{union}}$ 
24: end procedure

```

---

**Algorithm 2** minHashSig Calculation

---

```

1: procedure MINHASHSIG(termsBinaryRepresentation)
2:   Create minHash array in size of numPermutations
3:   for File  $f \leftarrow \text{termsBinaryRepresentation.KeySet()}$  do
4:     bins  $\leftarrow \text{termsBinaryRepresentation.get}(f)$ 
5:     for Permutation  $p \leftarrow \text{permutationses}$  do
6:       min  $\leftarrow \text{Integer.MaxValue}$ 
7:       for Boolean  $b \leftarrow \text{bins}$  do
8:         if  $b$  is true then
9:           if  $p.(b) < \text{min}$  then
10:            min  $\leftarrow p.(b)$ 
11:          end if
12:        end if
13:      end for
14:    end for
15:    Store min in minHash
16:  end for
17:  return minHash
18: end procedure

```

---

---

**Algorithm 3** Approximate Jaccard Calculation

---

```

1: procedure APPROXIMAT JACCARD(file1, file2)
2:    $b1 \leftarrow \text{minHash of file1}$ 
3:    $b2 \leftarrow \text{minHash of file2}$ 
4:   intersection  $\leftarrow 0$ 
5:   for each pair  $\langle d1, d2 \rangle$  in  $b1 \ b2$  do
6:     if  $d1$  is true and  $d2$  is true then
7:       intersection  $\leftarrow$  intersection + 1
8:     end if
9:   end for
10:  return  $\frac{\text{intersection}}{\text{numPermutations}}$ 
11: end procedure

```

---



---

**Algorithm 4** MinHash Matrix Calculation

---

```

1: procedure MINHASH MATRIX(minHash)
2:   matrix  $\leftarrow [\text{numPermutations}][\text{min.Hash.Keyset.size}]$ 
3:   Iterator  $\leftarrow \text{minHash's iterator}$ 
4:   for each string in Iterator do
5:     hashes  $\leftarrow$  iterator.next
6:   foreach permutation  $e$  in permutationses
7:     matrix  $\leftarrow$  hashes[numPermutation]
8:   end for
9:   return matrix
10: end procedure

```

---

## MinHashAccuracy

Table 1: Error is row, Number of Permutation is column. Data is (Large Error Pairs, Exact Jac in Sec, Appro Jac in Sec)

	<b>0.04</b>	<b>0.07</b>	<b>0.09</b>
<b>400</b>	(3131, 0.527, 0.199)	(8, 0.535, 0.184)	(1, 0.565, 0.189)
<b>600</b>	(658, 0.52, 0.268)	(2, 0.535, 0.268)	(1, 0.563, 0.279)
<b>800</b>	(267, 0.532, 0.361)	(4, 0.542, 0.358)	(1, 0.573, 0.368)

In conclusion, Increasing number of permutation result in less error, however, it does not bring a large increasing of time complexity. It is because the increasing number is small to the number of terms. For example, It is no large time complexity different between calculation on 400 rows and 600 rows by computer.

Another point is that, the large error pairs are in the range of  $[0.04, 0.07]$ , since when we adjust the allowed error to 0.09, there is only 1 pair in that range. The error happens because of the permutation we picked may not be able to totally permute the binary representation.