

Document similarity detection using hasing

Carlos Bergillos, Antoni Rambla, Adrià Cabeza
Departament de Computació

December 13, 2018

Abstract

Our goal is to identify similarities between documents. We have used *Local-Sensitive Hashing* algorithms and a *k-shingles* and *minhash signatures* representation of documents to evaluate the effectivity of the similarity computed. Once determined the best parameters for the collection, we will give an algorithm to indentify the more similar documents.

Contents

1	Introduction	3
2	Theory	3
2.1	Concept of similarity	3
2.2	Representation of documents	4
2.2.1	k -Shingles	4
2.2.2	Hashing Shingles	5
2.3	MinHash	5
2.4	Locality-Sensitive Hashing for Documents	6
2.4.1	Hash function used to hash a vector	6
2.4.2	Filling the signature matrix	6
2.4.3	Get the similarity	6
2.5	Distance Measures	8
2.5.1	Jaccard Distance	8
2.5.2	Edit Distance	8
2.6	Hashing	8
2.6.1	Modular Hashing	9
2.6.2	Multiplicative Hashing	9
3	Data	10
3.1	Real-world data	10
3.2	Generating the data	10
4	Experiments	11
5	Conclusion	11

1 Introduction

Our goal is to identify similarities between documents. We say that two documents are similar if they contain a significant number of common substrings that are not too small.

The problem of computing the similarity between two files has been studied extensively and many programs have been developed to solve it. Algorithms for the problem have numerous applications, including spelling correction systems, file comparison tools or even the study of genetic evolution.

Existing approaches can also include a brute force approach of comparing all sub-strings of pair of documents. However, such an approach is computationally prohibitive.

2 Theory

He cometido el peor de los pecados que un hombre puede cometer. No he sido feliz. Que los glaciares del olvido me arrastren y me pierdan, despiadados. Mis padres me engendraron para el juego arriesgado y hermoso de la vida, para la tierra, el agua, el aire, el fuego.

2.1 Concept of similarity

First we have to focus into the definition of similarity, when we talk about the “Jaccard similarity”, which is calculated by looking at the relative size of their intersection.

The Jaccard similarity, also known as Jaccard index is a statistical measure of similarity of sets. For two sets, it is defined as the size of the intersection divided by the size of the union of the sample sets. Mathematically,

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

Calculating the similarity estimation using this approach could be solved using k independent repetitions of the MinHash algorithm, however this would require $O(k \cdot |A|)$ running time.

The complexity of our implementation is $O(n)$. To calculate the intersection of both sets we use two iterators that iterate through both sets. When an element of a set is smaller than the other we increment its iterator and when they are equal we iterate both and a counter. Finally we use the value of the counter to apply the formula.

```
double intersection(const set<string>& A, const set<string>& B){
    int common = 0;
    auto it = A.begin(), it2 = B.begin();
    while(it != A.end() and it2 != B.end()){
        if(*it < *it2){
            ++it;
        }
```

```

    }
    else if(*it > *it2){
        ++it2;
    }
    else {
        ++it;
        ++it2;
        ++common;
    }
}
return common;
}

double Jaccard(const set<string>& A, const set<string>& B){
    double size_intersection = intersection(A,B);
    return size_intersection / (A.size() + B.size() -
        ↪ size_intersection);
}

```

If we take in account the cost of building a set we should observe that for unsorted sequences our cost would be incremented to $O(n * \log(N))$.

2.2 Representation of documents

To identify lexically similar documents we need a proper way to represent documents as sets and the most effective way is to construct from the document the set of short strings that appear within it. If we do so, even if the documents have different sizes or those sentences appear in different order we will find several common elements. In the next section we will introduce some of the approaches of shingling and its variations.

2.2.1 *k*-Shingles

A *k*-shingle (or word-*k*-gram) is a sequence of consecutive words of size *k*. Intuitively, two documents A and B are similar if they share enough *k*-shingles. By performing union and intersection operations between the *k*-shingles, we can find the Jaccard similarity coefficient between A and B.

There are some variations regarding on how white space (blank, tab, newline, etc) is treated. Also there is a variation that works with a bag of shingles instead of a set to keep the number of appearances of a shingle.

How large *k* should be depends on how long typical documents are and how large the set of typical characters is. For example if we pick *k*=4 there are $27^4 = 531441$ possible

k-shingles. However, the calculation can be a little bit more subtle because all the characters do not appear with equal probability. A good rule of thumb is to imagine that there are only 20 characters and estimate the number of k-shingles as 20^k . For large documents, choice $k = 9$ is considered safe.

2.2.2 Hashing Shingles

Instead of using substrings directly as shingles, we can pick a hash function that maps strings of length k to some number of buckets and treat the resulting bucket number as the shingle. That process compacts our data and lets us manipulate shingles by single-word machine operations.

2.3 MinHash

A minhash function on sets is based on a permutation of the universal set. Given any such permutation, the minhash value for a set is that element of the set that appears first in the permuted order.

This algorithm provides us with a fast approximation to the Jaccard Similarity. The concept is to condense the large sets of unique shingles into a much smaller representations called “signatures”. We will then use these signatures to measure the similarity between document, the signature won’t give us the exact similarity but we will get a close estimate (the larger the number of signatures you choose, the more accurate the estimate).

To implement the idea of generating randomly permuted rows, we don’t actually generate the random numbers, since it is not feasible to do so for large datasets, e.g. For a million itemset you will have to generate a million integers . . . , not to mention you have to do this for each signatures that you wish to generate. One way to avoid having to generate n permuted rows is to pick n hash functions in the form of :

$$h(x) = (ax + b) \bmod(c)$$

Where:

- x is the row numbers of your original characteristic matrix
- a and b are any random numbers smaller or equivalent to the maximum number of x
- c is a prime number slightly larger than the total number of shingle sets.

2.4 Locality-Sensitive Hashing for Documents

This technique allows us to avoid computing the similarity of every pair of sets or their minhash signatures. If we are given signatures for the sets, we may divide them into bands, and only measure the similarity of a pair of sets if they are identical in at least one band. By choosing the size of bands appropriately, we can eliminate from consideration most of the pairs that do not meet our threshold of similarity.

2.4.1 Hash function used to hash a vector

```
int hash_vec(V) {
    seed = V.size();
    for each i in V do:
        seed = seed XOR (i + 0x9e3779b9 + (seed << 6) + (seed >> 2));
        //operands << and >> are shifting left and right respectively
        ↪ x bits
    return seed;
}
```

2.4.2 Filling the signature matrix

```
fill(Matrix<unsigned int> repMatrix, set<string> shingles,
vector<set<string>> docShing){
    set iterator it = shingles.begin();
    for i = 0 to repMatrix.size() do:
        repMatrix[i][0] = i;
        string shingle = shingles.get(it);
        it++;
        for j = 1 to repMatrix[0].size() do:
            if(docShing[j-1].contains(shingle))then repMatrix[i][j] =
                ↪ 1;
            //el programa en columna j conte el shingle "shingle"
            else repMatrix[i][j] = 0;
    }
```

2.4.3 Get the similarity

```
float sim(Matrix<unsigned int> signatureMatrix, int a, int b)
    float simil = 0;
    for int i = 0 to signatureMatrix.size() do:
        if(signatureMatrix[i][a] == signatureMatrix[i][b])
            then ++simil;
```

```
    return simil / signatureMatrix.size();
```

Funció que ens ordena un parell posant el més gran com a segon.

```
parella_inc(int a, int b){  
    if(a < b ) then return pair(a,b);  
    else return pair(b,a);
```

```
main()  
    set<string> shingles ,aux;  
    //set of the total of shingles we have  
    and an auxiliar for each document  
    vector<set<string>> docShing;  
    //for each document we save its set of shingles  
    for i = 1 to number of documents do:  
        aux = kshingles(file[i], 9, false , true , true);  
        shingles.add(aux);  
        docShing[i] = aux;  
    Matrix<unsigned int> repMatrix (shingles.size() ,  
        vector<unsigned int> (number of files +1 ));  
    fill(repMatrix , shingles , docShing);  
    int b, r, h;  
    cin >> b >> r; //number of bands and rows  
    h = r*b;  
    Matrix<unsigned int> signatureMatrix  
        (h, vector<unsigned int> (number of files , INFINITY));  
    hashing of choice  
    set<pair<unsigned int ,unsigned int>> candidats;  
    map<unsigned int ,vector<unsigned int>> bucket;  
    for i = 0 to h do:  
        bucket.clear();  
        for j = 0 to number of documents do:  
            vector<unsigned int> row;  
            for k = 0 to r do:  
                row.push_back(signatureMatrix[i+k][j]);  
                //computing rows  
            unsigned int doc1 = hash_vec(row);  
            if(bucket.containsKey(doc1)) then  
                //we find a row that has the same hash signature  
                for l = 0 to bucket[doc1].size() do:  
                    //we may have several rows with the same hash signature
```

```

    ↪ that
    //means they are all candidates with each other
    candidats.insert(parella_inc(bucket[doc1][1],j));

    bucket[doc1].push_back(j);
    else
        bucket[doc1] = vector<unsigned int>(1, j);
}

```

2.5 Distance Measures

In this section we shall take a break to study the notions of distance measure. Suppose we have two points (i.e. x and y), a distance measure on them is a function $d(x,y)$ that produces a real number and satisfies the following properties:

- $d(x,y) \geq 0$
- $d(x,y) = 0 \iff x = y$
- $d(x,y) \leq d(x,z) + d(z,y)$ (this is called the triangle inequality)

2.5.1 Jaccard Distance

We define the Jaccard distance of two sets x and y as $d(x,y) = 1 - \text{SIM}(x,y)$. Note that the Jaccard similarity is the probability that a random hash function maps x and y to the *same* value, that means $d(x,y)$ is the probability that a random hash function sends them to *different* values.

2.5.2 Edit Distance

2.6 Hashing

Hashing is a technique for dimensionality reduction. It uses a hash function that is any function that can be used to map data (called a key) of arbitrary size to data of a fixed size (called a hash value or hash). That hash is a sum up of everything that is in the data. You can never make it backwards from the hash to the data.

Hashing is done for indexing and locating items in databases because it is easier to find the shorter hash value than the longer string.

The hash functions that we want to use need to be:

- **Really fast**
Dimensionality reduction is often a time bottle-neck and using a fast basic hash function to implement it may improve running times significantly.

- **Avoid hash collisions**

We do not want to get the same hash using different pieces of data

- **Uniform distribution**

The hash values are uniformly distributed.

In our project we have not chosen hash cryptographic functions (i.e. sha-1 or md5) because they are too slow for our purpose.

2.6.1 Modular Hashing

In the modular hashing (also called the division method), we map a key k into one of m slots by taking the remainder of k divided by m . It takes the following form.

$$h(k) = k \bmod(m)$$

When using the division method, we usually avoid certain values of m . For example, $m = 2^p$ for some integer p , then $h(k)$ would be just the p -lowest-order bits of k . A prime value is often a good choice of m .

2.6.2 Multiplicative Hashing

The multiplicative method for creating hash functions operates in two steps. Firstly, we multiply the key k by a constant A in the range $0 < A < 1$ and extract the fractional part of kA . Then, we multiply this value by m and take the floor the result. To sum up:

$$\lfloor m \cdot kA \bmod(1) \rfloor$$

Where:

- $kA \bmod(1)$ is the fractional part of kA , that is, $kA - \lfloor kA \rfloor$

An advantage of the multiplication method is that the value of m is not critical. We typically choose it to be a power of 2 ($m = 2^p$ for some integer p), since we can then easily implement the function on most computers. Supposing that the word size of the machine is w bits and that k fits in a single word. We restrict A to be a fraction of the form $S/2^w$, where s is an integer in the range $0 < s < 2^w$. We first multiply k by the w -bit integer $s = A \cdot 2^w$. The result is a $2w$ -bit value $r_1 2^w + r_0$, where r_1 is the high-order word of the product and r_0 is the lower-order word of the product. The desired p -bit hash value consists of the p most significant bits of r_0 .

FICAR FIGURA CARLOS AMB ILLUSTRATOR

FICAR EL PSEUDOCODI

Although this method works with any value of the constant A , it works better with some values than with others. The optimal choice depends on the characteristics of the data being hashed. Knuth suggest that...

$$A = (\sqrt{5} - 1)/2 \simeq 0.6180339887$$

...is likely to work well.

3 Data

3.1 Real-world data

- **Harry Potter and the Sourcerer Stone:** All the text from the first Harry Potter novel.
- **The Lord of The Rings: The return of the King:** The entire script from the last Lord of the Rings movie.
- **Star Wars: Heir to the Empire:** An Star Wars novel by Timothy Zahn.

3.2 Generating the data

To generate the data for our experiments we have used the `random_shuffle` function inside the C++ STL which rearranges the elements randomly. The function swaps the value of each element with that of some other randomly picked element. Its complexity is $O(n)$.

```
ifstream inFile;
ofstream outFile;

void writeVector(const vector<string>& A, int i, string nameFile){
    outFile.open(nameFile);
    for(int j = 0; j < A.size(); ++j){
        outFile << A[j]<< space symbol;
    }
    outFile.close();
}

int main(){
    int n;
    cin>>n; //number of permutations that we want to create

    string nameFile;
    cin>> nameFile;
```

```

inFile.open(nameFile); //open the file

vector<string> words;
string testline;
while( inFile >> testline){
    words.push_back(testline);
    //insert each word in a vector of strings
}
inFile.close();
for(int i = 0; i<n;++i){
    random_shuffle(words.begin(),words.end());
    writeVector(words,i,nom);
}
}

```

As it can be seen, we traverse all the document once to generate a vector of strings containing all the words and then we shuffle the vector and write the result in a document n times.

4 Experiments

Vamos a poner unos cuantos experimentos

5 Conclusion

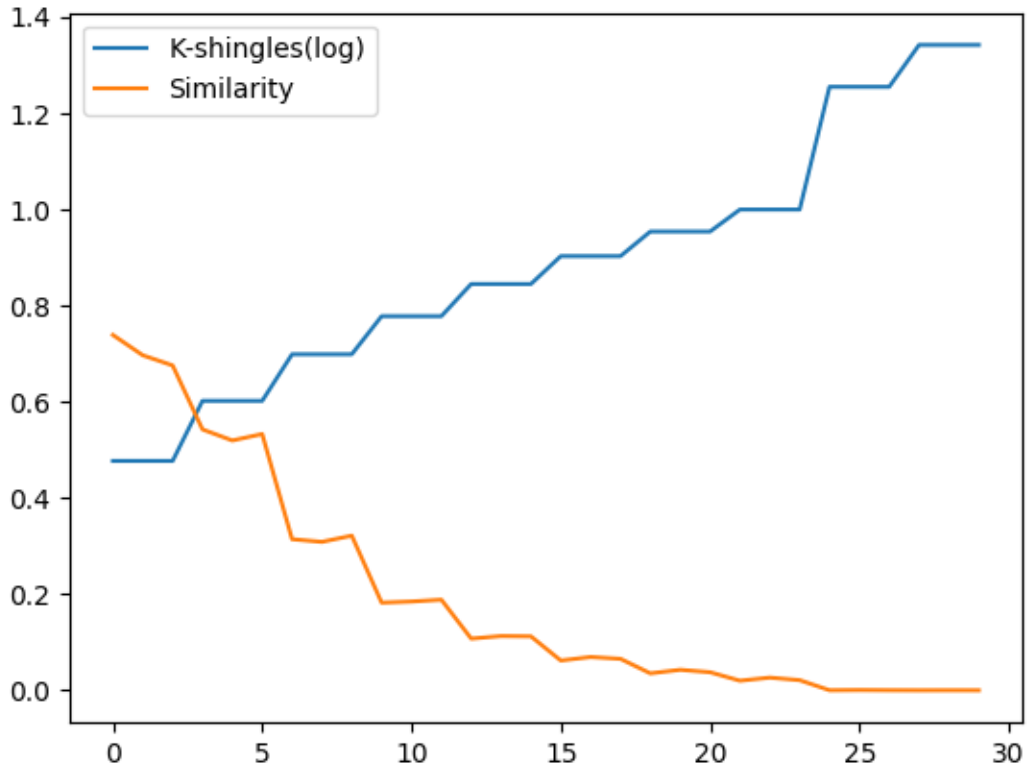


Figure 1: Relation between the k-shingles and the similarity outputted

References

- [1] J. Bank and B. Cole, *Calculating the Jaccard similarity coefficient with map reduce for entity pairs in Wikipedia* (Wikipedia Similarity Team, 2008).
- [2] J. Leskovec, A. Rajaraman and J. Ullman *Finding Similar Items. In Mining of Massive Datasets* (Cambridge University Press, 2014).
- [3] E. W. Myers *An $O(ND)$ Difference Algorithm and Its Variations* (Department of Computer Science, University of Arizona)
- [4] S. Alzahrani and N. Salim, *Fuzzy Semantic-Based String Similarity for Extrinsic Plagiarism Detection* (Taif University, Saudi Arabia and Universiti Teknologi Malaysia, Malaysia)

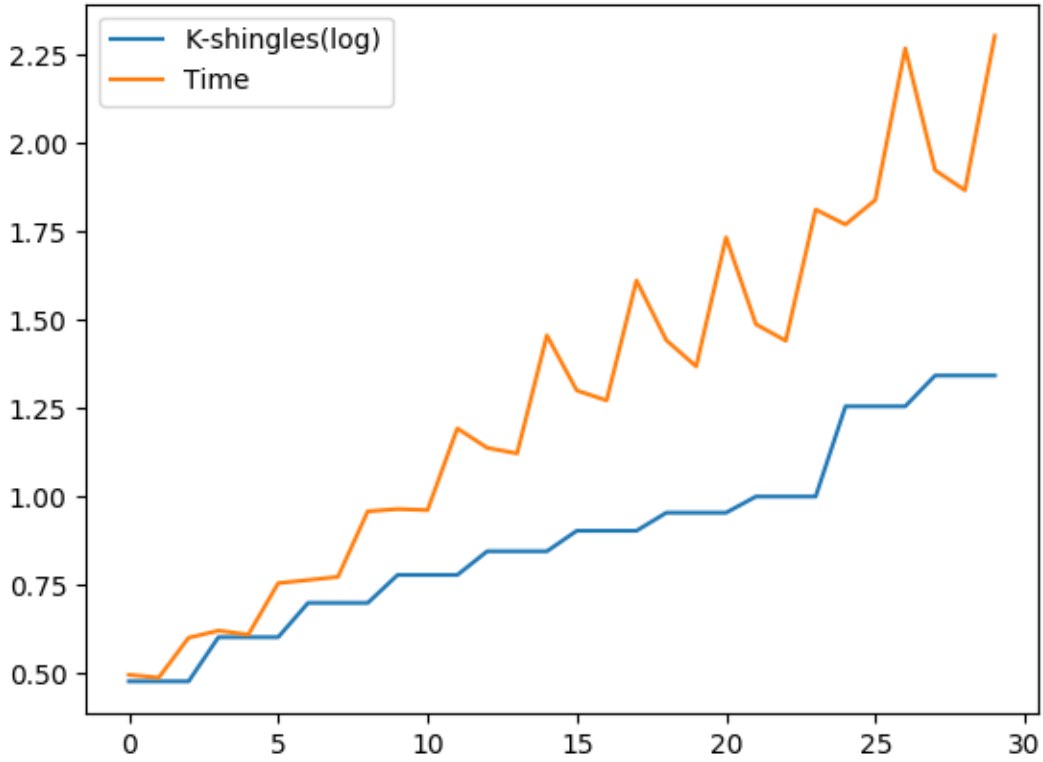


Figure 2: Relation between the k-shingles and the time of execution

- [5] S. Dahlgaard, M. Tejs Knudsen and M. Thorup, *Practical Hash Functions for Similarity Estimation and Dimensionality Reduction* (University of Copenhagen)
- [6] D. Lemire and O. Kaser, *Strongly universal string hashing is fast* (Université du Québec Canada and University of New Brunswick, Canada)
- [7] C++ Standard Template Library, Source: cplusplus.com/reference/stl/
- [8] T. H. Cormen, C. E. Leieron, R. L. Rivest and C. Stein, *Introduction to Algorithms. Third Edition*
- [9] D. Knuth, *Sorting and Searching*, volume 3 of *The Art of Computer Programming* (Addison-Wesley, 1997. Third Edition)