



# Teoría de Algoritmos 2

Segundo Cuatrimestre 2017

## Trabajo Práctico 1

Integrante	Padrón	Correo electrónico
Rodrigo De Rosa	97799	rodrigoderosa@outlook.com
Marcos Schapira	97934	schapiramarcos@gmail.com
Facundo Guerrero	97981	facundoiguerrero@gmail.com

# Índice

<b>1. Rabin - Karp</b>	<b>1</b>
1.1. Funcionamiento . . . . .	1
1.2. Implementación . . . . .	1
1.3. Complejidad . . . . .	1
1.4. Investigación y aplicaciones . . . . .	1
1.5. Conclusiones . . . . .	1
<b>2. Zhu-Takaoka</b>	<b>1</b>
2.1. Funcionamiento . . . . .	1
2.2. Implementación . . . . .	1
2.3. Complejidad . . . . .	2
2.4. Investigación y aplicaciones . . . . .	2
2.5. Conclusiones . . . . .	2

# 1. Rabin - Karp

*Algoritmo Michael O. Rabin and Richard M. Karp 1987*

## 1.1. Funcionamiento

La idea del algoritmo es muy simple. Basándose en la estructura del algoritmo naïve, este agrega un paso previo que compara los strings por valores de hash. Para esto precisa una función de hash que se busca que compare entre valores lo mas rápido posible. Esto tiene el potencial beneficio de acortar los tiempos de comparación entre strings mientras que agrega la complejidad del calculo previo del valor de hash para cada string.

## 1.2. Implementación

La implementación es muy simple. Primero calcula el valor de hash para el patrón a buscar. Luego recorre el texto calculando el valor de hash para la palabra a buscar. Compara ambos valores y si dan iguales entonces compara si las palabras son realmente iguales o no. Para ganar mayor velocidad se utilizo la librería pyhash <sup>1</sup> que contiene implementaciones en C/C++ para mejor eficiencia de algoritmos no criptográficos. De estos se usaron (todos de 32 bits): FNV, Murmur Hash, City Hash, Spooky Hash.

## 1.3. Complejidad

En el peor de los casos, el algoritmo compara cada string del texto contra el patrón teniendo un orden de  $O(nm)$  donde  $n$  es la longitud del texto y  $p$  la del patron. Esto ocurre en el caso en donde se use una función de hash muy mala. Con una función de hash relativamente buena se mejora el orden a  $O(n + m)$ .

## 1.4. Investigación y aplicaciones

Este algoritmo no es utilizado para Simple Matching ya que resulta poco eficiente. Esto se debe que el costo que tiene para calcular las claves entre algoritmos resulta mayor en relación al beneficio que se obtiene de la rapidez para comparar strings. Investigando sobre sus aplicaciones en el ámbito profesional se encuentra que este algoritmo resulta particularmente útil para el problema de múltiple string matching, mas es así en la búsqueda de plagios. Esto es, teniendo un texto A y un texto B, comparar que tan semejante resulta A contra B.

## 1.5. Conclusiones

Para simple matching este algoritmo resulta increíblemente ineficiente dando los peores tiempos ejecución. Sin embargo para múltiple matching es un muy buen algoritmo. Como optimización se sigue sacando la parte en donde se verifica que los valores de hashes que tuvieron match sean realmente iguales. Esto funcionaría sin problemas con una función de hashing perfecta (pero al entiza la ejecución), sin embargo si no lo es el algoritmo pasaría a ser randomizado ya que las funciones de hash utilizadas en este caso garantizan pocas colisiones pero no es imposible que ocurran.

# 2. Zhu-Takaoka

*Algoritmo Zhu Rui Feng - Tadao Takaoka 1987*

## 2.1. Funcionamiento

El algoritmo que esta siendo presentado es una variante del algoritmo de Boyer-Moore. Este algoritmo, al igual que el de BM, mantiene la regla de “good suffix” pero reemplaza la regla de “bad character” por la regla de “2-substrings”. Lo que hace esta ultima regla es guardar en una matriz las apariciones mas a la derecha de cada par de caracteres (a,b) pertenecientes al patrón. Entonces el algoritmo va a comparar el texto con el patrón aplicando una de las 2 reglas en caso de encontrar un miss o un match.

## 2.2. Implementación

El algoritmo consta de 2 fases. La primera fase es la de pre-procesamiento en donde se calcula la matriz necesaria para aplicar la regla “2-substrings” y donde se crea el vector para la regla “good suffix” al igual que en el algoritmo de Boyer-Moore. En la segunda fase, el algoritmo alinea el texto con el patrón a izquierda y recorre de derecha a izquierda el patrón comparando carácter a carácter con el texto. En caso de encontrar un miss o de llegar a un

---

<sup>1</sup><https://github.com/flier/pyfasthash>

match, el algoritmo calcula el máximo entre las 2 reglas antes mencionadas, y shiftea el patrón a derecha en esa cantidad. Esto se repite, hasta que el patrón llega al final del texto.

### 2.3. Complejidad

Este algoritmo tiene una complejidad de  $O(m + a^2)$  para tiempo y espacio en la fase de pre-procesamiento, siendo  $a$  el tamaño del alfabeto. Pero para la fase de búsqueda el algoritmo tiene una complejidad temporal de  $O(nm)$ , siendo  $n$  y  $m$  el tamaño del patrón y del texto respectivamente.

### 2.4. Investigación y aplicaciones

Este algoritmo es utilizado con alfabetos pequeños, ya que es cuando resulta eficiente. Esto es debido a la dependencia del tamaño del alfabeto con la fase de pre-procesamiento. Además, este algoritmo resulta muy eficiente en multiple string matching en 2 dimensiones.

### 2.5. Conclusiones

El algoritmo anteriormente presentado, es uno de los que mejor tiempos tiene dentro de los algoritmos implementados en dicho trabajo. Además, se puede ver claramente que la fase de pre-procesamiento aumenta abruptamente a medida que aumentamos el tamaño del alfabeto, tanto en espacio como en tiempo. Se concluye, que este algoritmo funciona muy rápidamente cuando el alfabeto o el patrón son chicos. Como recomendación adicional, se aconseja utilizarlo para patrones chicos o multiple matching.