



INSTITUTO
DE INGENIERÍA
UNAM

Métodos para la obtención automática de términos en un área de especialidad



Ing. Alberto Barrón Cedeño

Instituto de Ingeniería

Posgrado en Ciencia e Ingeniería de la Computación

UNAM

alberto@pumas.ii.unam.mx



- Introducción
- Técnicas para la extracción de términos
- El algoritmo C-value/NC-value
 - Parte lingüística de C-value
 - Parte estadística de C-value
 - NC-value
- Un breve ejemplo
- Conclusiones



Proyecto CONACyT:

Extracción de conceptos en textos de especialidad a través del reconocimiento de patrones lingüísticos y metalingüísticos

Creación de diccionarios

Creación de tesauros y ontologías

Desarrollo de buscadores

Clasificadores de documentos



Técnicas lingüísticas

- Los términos siguen ciertos patrones sintácticos.
- Se requiere etiquetado POS.

Ejemplos: Lexter, Heid



Técnicas estadísticas

- Un suceso que ocurre frecuentemente es relevante

Ejemplo: ANA



- Técnicas híbridas
 - Combinan las dos técnicas anteriores

Ejemplos: Acabit, TermoStat, Termine



Es un método híbrido para la extracción automática de términos multipalabra en inglés.

En particular, se ha desarrollado para el procesamiento de documentos del área de biomedicina.



C-value (única etapa híbrida)

lingüística

Etiquetado de partes de la oración

Detección de candidatos a término

Eliminación de candidatos por medio de una lista de paro

estadística

Ordenación de candidatos con base en su frecuencia y longitud



Texto original

Tal y como sale de fábrica el disco duro no puede ser utilizado por un sistema operativo.

Antes tenemos que definir en él una o más particiones y luego hemos de darles un formato que pueda ser entendido por nuestro sistema.



Texto etiquetado

Tal tal/QU y y/CC como como/CSUBX sale_salar/VLfin
de_de/PDE fábrica fábrica/NC el el/ART disco disco/NC
duro duro/ADJ no no/NEG puede_poder/VMfin ser_ser/VSinf
utilizado_utilizar/VLadj por_por/PREP un_un/ART
sistema_sistema/NC operativo_operativo/ADJ ._/FS

Antes antes/ADV tenemos tener/VLfin que_que/CQUE
definir_definir/VLinf en_en/PREP él_él/PPX una_un/ART o_o/CC
más_más/ADV particiones_particiones/NC y_y/CC
luego_luego/CSUBF hemos_haber/VLfin de_de/PDE
darles_dar/VLinf un_un/ART formato_formato/NC que_que/CQUE
pueda_poder/VMfin ser_ser/VSinf entendido_entender/VLadj
por_por/PREP nuestro_nuestro/PPO sistema_sistema/NC ._/FS



$\langle \text{NC} \mid \text{NP} \mid \text{PE} \rangle^+$

servidor, tarjeta madre, MB, arquitectura von Neumann

$\langle \text{NC} \rangle \langle \text{ADJ} \rangle (\langle \text{PDE} \rangle \langle \text{NC} \mid \text{NP} \rangle)^*$

sistema operativo, unidad central de procesamiento

$\langle \text{NC} \rangle \langle \text{PDE} \rangle \langle \text{NC} \mid \text{NP} \mid \text{NMEA} \rangle$

ángulo de visión, ancho de banda, tasa de transferencia



<VLFIN | VLINF>

compilar, descifrar, ensamblar

<NC>? <ACRNM>

DDR2, IP, slot ISA

<NC> <PDE> ((<NC> <ADJ>) | (<ADJ> <NC>))

computadora de uso general



Tal y como sale de **fábrica** el **disco duro** no puede ser utilizado por un **sistema operativo**.

Antes tenemos que definir una o más **particiones** en el **disco duro** y **formatearlo** para que pueda ser entendido por nuestro **sistema**.



Más de 200 palabras (nombres y adjetivos) que no se espera que aparezcan dentro de los términos del área.

Ej. (computación): detalle, elegante, importante, mayoría, mercado, opción, tamaño, vez



~~fábrica~~

disco duro

sistema operativo

partición

formato

sistema





Aspectos considerados

1. Frecuencia total de ocurrencia del sintagma en el corpus
2. Frecuencia total de ocurrencia del sintagma como parte de sintagmas más largos
3. Número de dichos candidatos de mayor longitud
4. Longitud del candidato a término



$$C - value = \begin{cases} (1 + \log_2 |a|) * f(a) & \text{si } a \text{ no aparece en otros candidatos} \\ (1 + \log_2 |a|) * \left(f(a) - \frac{1}{P(T_a)} \sum_{b \in T_a} f(b) \right) & \text{de lo contrario} \end{cases}$$

a es el sintagma candidato

$|a|$ es la longitud de a

$f(a)$ es la frecuencia de ocurrencia de a en el corpus

T_a es el conjunto de candidatos de mayor longitud que contienen a a

$P(T_a)$ es el número de esos candidatos (incluye al mismo candidato)

$\sum f(b)$ es la ocurrencia total de a como subcadena del sintagma candidato b tal que $|a| < |b|$



candidato	longitud	frecuencia
disco duro	2	3
sistema operativo	2	2
sistema	1	2
partición	1	1



candidato	longitud	frecuencia
sistema operativo	2	2

$$C - value = (1 + \log_2|a|) * f(a)$$

$$C - value(sistema operativo) = (1 + \log_2(2)) * 2$$

$$C - value(sistema operativo) = 4$$



candidato	longitud	frecuencia
sistema	1	2

$$C - value = (1 + \log_2|a|) * \left(f(a) - \frac{1}{P(T_a)} \sum_{b \in T_a} f(b) \right)$$

$$C - value(sistema) = (1 + \log_2(1)) * \left(2 - \frac{1}{2} * 2 \right)$$

$$C - value(sistema) = 2$$



candidato	longitud	frecuencia
disco duro	2	3
sistema operativo	2	2
sistema	1	2
partición	1	1



candidato	longitud	frecuencia	C-value
disco duro	2	3	6
sistema operativo	2	2	4
partición	1	1	1
sistema	1	2	1



El contexto en el que se hallan los candidatos es importante para ellos.

Las palabras que actúan con los términos no suelen ser arbitrarias.

Son **palabras de contexto** aquellos nombres, adjetivos y verbos que aparecen en el contexto de un candidato.



Un **llavero USB** es un **pequeño dispositivo** de almacenamiento que utiliza la memoria flash para guardar la información sin necesidad de pilas.

Tienen una capacidad de almacenamiento que va desde algunos megabytes hasta 8 gigabytes.



Un **llavero USB** es un pequeño **dispositivo de almacenamiento** que **utiliza la memoria flash** para guardar la información sin necesidad de pilas.

Tienen una capacidad de almacenamiento que va desde algunos megabytes hasta 8 gigabytes.



Un llavero USB es un pequeño dispositivo de **almacenamiento** que **utiliza** la **memoria flash** para **guardar** la **información** sin necesidad de pilas.

Tienen una capacidad de almacenamiento que va desde algunos megabytes hasta 8 gigabytes.



Un llavero USB es un pequeño dispositivo de almacenamiento que utiliza la memoria flash para guardar la información sin necesidad de pilas.

Tienen una **capacidad de almacenamiento** que va desde algunos **megabytes** hasta 8 gigabytes.



Un llavero USB es un pequeño dispositivo de almacenamiento que utiliza la memoria flash para guardar la información sin necesidad de pilas.

Tienen una capacidad de **almacenamiento** que va desde algunos **megabytes** hasta 8 **gigabytes**.



Un llavero USB es un pequeño dispositivo de almacenamiento que utiliza la memoria flash para guardar la información sin necesidad de pilas.

Tienen una capacidad de almacenamiento que va desde algunos **megabytes** hasta 8 **gigabytes**.



candidato

disco duro

sistema operativo

partición

sistema

palabras de contexto

formatear, guardar, GB

instalar, configurar, cargar

disco duro, crear, swap

levantar, iniciar, apagar



$$\textit{weight}(w) = \frac{t(w)}{n}$$

w es la palabra de contexto analizada

$\textit{weight}(w)$ es el peso asignado a la palabra w

$t(w)$ es el número de candidatos con los que aparece la palabra w

n es el número total de candidatos considerados (para expresarlo como una probabilidad)



“es necesario **formatear** el *disco* y *crear* dos **particiones**”

“el tipo de **partición** para el *sistema operativo* **Linux es ext3** ”

$$\textit{weight}(\textit{formatear}) = \frac{1}{2} \quad \textit{weight}(\textit{particion}) = 1$$



Tal y como sale de **fábrica** el **disco duro** no puede ser utilizado por un **sistema operativo**.

Antes tenemos que definir una o más **particiones** en el **disco duro** y formatearlo para que pueda ser entendido por nuestro **sistema**.



Tal y como **sale** de **fábrica** el **disco duro** no **puede ser utilizado** por un **sistema operativo**.

Antes tenemos que definir una o más **particiones** en el **disco duro** y formatearlo para que pueda ser entendido por nuestro **sistema**.



$$NC - value(a) = 0.8C - value(a) + 0.2 \sum_{b \in C_a} f_a(b)weight(b)$$

$$NC - value(disco duro) = 0.8(6) + 0.2 \left(1 * 1 + 1 * \frac{1}{2} \right)$$

$$NC - value(disco duro) = 5$$



candidato	frecuencia	C-value	NC-value
disco duro	3	6	5.3
sistema operativo	2	4	5.0
formato	2	2	1.8
partición	1	1	1.4
sistema	2	1	0.8



candidato	frecuencia	C-value	NC-value
usuario	1	1	1
estación de trabajo	69	7	8
problema	119	140	138
memoria flash	176	85	82



USB El Bus de Serie Universal (USB, de sus siglas en inglés Universal Serial Bus) es un interfaz que provee un estándar de bus serie para conectar dispositivos a un ordenador personal (generalmente a un PC). Un sistema USB tiene un diseño asimétrico, que consiste en un solo servidor y múltiples dispositivos conectados en una estructura de árbol utilizando concentradores especiales. Se pueden conectar hasta 127 dispositivos a un solo servidor, pero la suma debe incluir a los concentradores también, así que el total de dispositivos realmente usables es algo menor. Fue desarrollado a finales de 1996 por siete empresas: IBM, Intel, Northern Telecom, Compaq, Microsoft, Digital Equipment Corporation y NEC. El estándar incluye la transmisión de energía



El algoritmo C-value/NC-value ha mostrado ser una buena opción en búsqueda de la extracción de términos en el español.

Los errores pueden ser “soportados” si se considera que una persona puede tardar unas 3 horas en obtener la terminología de un documento de 2381 palabras.

Se tarda unas horas en obtener la terminología de un corpus de 140,000.





**INSTITUTO
DE INGENIERÍA
UNAM**

¡Gracias!



Alberto Barrón Cedeño
alberto@pumas.ii.unam.mx

3er. Coloquio de Lingüística Computacional en la UNAM