

DESCUBRIMIENTO DE RELACIONES PARADIGMÁTICAS

Salvador López Mendoza

Mayo de 2018

CONTEXTO DE PALABRAS

Se define un *seudo-documento*.

Al considerar las siguientes oraciones:

Mi gato come pez los sábados.

Su gato come pavo los martes.

Se tiene:

$$Izq(gato) = \{mi, su, el, un, \dots\}$$

$$Der(gato) = \{come, es, tiene, \dots\}$$

$$Todo(gato) = \{mi, su, el, un, come, es, tiene, pez, pavo, martes, \dots\}$$

SEUDO DOCUMENTO

Contexto es un pseudo documento.

El pseudo documento se maneja como una *bolsa de palabras*.

En una bolsa las palabras no tienen orden.

El contexto puede contener palabras que son adyacentes, pero también palabras que no lo son.

SIMILARIDAD DE CONTEXTOS

Medir la similaridad de contextos:

$$\begin{aligned}\text{Sim}(\text{perro}, \text{gato}) &= \text{Sim}(\text{Izq}(\text{perro}), \text{Izq}(\text{gato})) \\ &\quad + \text{Sim}(\text{Der}(\text{perro}), \text{Der}(\text{gato})) \\ &\quad \dots \\ &\quad + \text{Sim}(\text{Todo}(\text{perro}), \text{Todo}(\text{gato}))\end{aligned}$$

Un valor alto en $\text{Sim}(\text{palabra}_1, \text{palabra}_2)$ quiere decir que palabra_1 y palabra_2 están *relacionadas paradigmáticamente*.

MODELO DE ESPACIO VECTORIAL (VSM)

La bolsa de palabras se puede representar como un *espacio vectorial*.

La bolsa de palabras correspondiente a *gato* es el conjunto de palabras que aparecen en las oraciones que contienen la palabra *gato*.

Sea N la cantidad de palabras distintas (vocabulario).

d_1 representa a la bolsa de palabras correspondiente a *gato*.

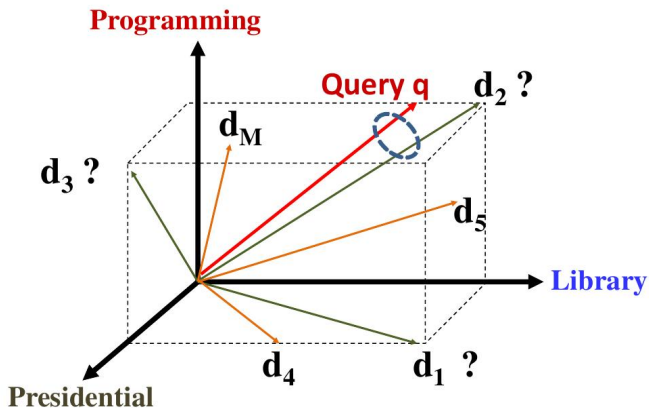
$$d_1 = (x_1, x_2, \dots, x_N)$$

Cada x_i indica la cantidad de veces que aparece esa palabra en la bolsa de palabras.

De la misma forma se define d_2 para la bolsa de palabras de *perro*.

$$d_2 = (y_1, y_2, \dots, y_N)$$

MODELO DE ESPACIO VECTORIAL (II)



SIMILARIDAD ENTRE DOCUMENTOS

¿Cómo se calcula la similaridad entre documentos?

$$\text{Sim}(d_1, d_2) = ?$$

Ejemplo:

Considerar **Su gato come pavo los martes**.

$$d_1 = (1, 1, 1, 1, 1)$$

¿Qué tan similar es a **Su perro come pavo los martes**?

$$d_2 = (1, 1, 1, 1, 1)$$

Comparar:

Su gato come pez los sábados con **Su perro come carne los domingos**.

Se necesita definir la función de similaridad.

¡Hay muchas posibilidades!

CALCULANDO LA SIMILARIDAD

$$\text{Sim}(d_1, d_2) = d_1 \cdot d_2$$

$$\text{Sim}(d_1, d_2) = x_1y_1 + x_2y_2 + \dots + x_ny_n$$

$$\text{Sim}(d_1, d_2) = \sum_{i=1}^N x_iy_i$$

Para el ejemplo anterior:

$$\text{Sim}(d_1, d_2) = 5$$

¿Cuál es la similaridad de las otras dos frases?

TRASLAPE ESPERADO DE PALABRAS EN EL CONTEXTO (EOWC)

Intuición: mientras mayor cantidad de palabras se traslapen, la similitud es mayor.

Problemas:

- Se favorece a aquellos elementos que aparezcan con mayor frecuencia.
- Se trata a todas las palabras equitativamente.
Las coincidencias en palabras como *los* no son tan importantes como las coincidencias en otras palabras, como *come*.

MEDIDAS DE SIMILITUD

Definiciones:

Sea $V = \{w_1, w_2, \dots, w_N\}$ el conjunto de palabras (*vocabulario*).

$c(w_i, d)$ es la cantidad de ocasiones en que aparece la palabra w_i en el documento d .

$$d_1 = (c(w_1, d_1), c(w_2, d_1), \dots, c(w_N, d_1))$$

Se normaliza el cálculo:

$$\text{Sea } x_i = c(w_i, d_1)/|d_1| \text{ y } y_i = c(w_i, d_2)/|d_2|$$

$$\text{Sim}(d_1, d_2) = d_1 \cdot d_2.$$

Notas:

x_i es la probabilidad de que al tomar al azar una palabra del documento d_1 , la palabra sea w_i .

$\text{Sim}(d_1, d_2)$ es la probabilidad de que al tomar al azar una palabra de cada documento, éstas sean idénticas.

RECUPERACIÓN DE TEXTOS (TR)

- Existe una colección de textos.
- El usuario hace una consulta para expresar lo que requiere.
- El sistema (*máquina de búsqueda*) regresa los documentos relevantes.

TR VS. BASES DE DATOS

- Información.
El texto no tiene estructura (es texto libre).
La información es ambigua, no hay una semántica bien definida.
- Consulta.
En TR es ambigua. En BD hay una semántica bien definida.
La especificación puede ser incompleta.
- Respuesta.
Documentos relevantes en TR. En BD son registros que coinciden.

DEFINICIÓN FORMAL DE TR

- Vocabulario. $V = \{w_1, w_2, \dots, w_N\}$
- Consulta. $q = q_1, \dots, q_m$, en donde $q_i \in V$
- Documento. $d_i = d_{i1}, \dots, d_{imi}$, en donde $d_{ij} \in V$
- Colección. $C = \{d_1, \dots, d_M\}$
- Conjunto de documentos relevantes. $R(q) \subset C$
Es un conjunto desconocido. Depende del usuario.
La consulta es una *sugerencia* de lo que debe contener un documento en $R(q)$.
- Tarea. Calcular $R'(q)$. Es una aproximación a $R(q)$.

Se usa una *función de clasificación* (ranking function).

Posibilidad: modelos basados en similitud.

SIMILITUD EN TR

- Vector de bits.

$$x_i, y_j \in \{0, 1\}$$

Al calcular $Sim(d_1, d_2)$ se cuenta la cantidad de palabras que coinciden.

- Vector de frecuencia de términos.

Se debe dar mayor valor al hecho de que una palabra se repita varias veces.

x_i = cantidad de ocasiones en que aparece w_i en la consulta.

y_i = cantidad de ocasiones en que aparece w_i en el documento.

Da más peso a las palabras que se repiten con frecuencia.

Agregar *frecuencia inversa*.

Ayuda a discriminar palabras que son frecuentes, pero no aportan contenido.

$$\text{Ahora } y_i = c(w_i, d) \times IDF(w_i)$$

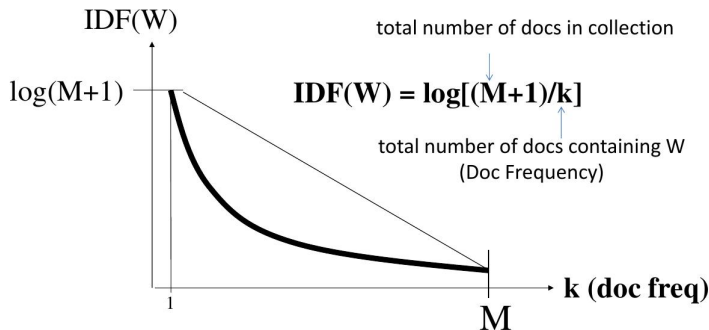
IDF

Penaliza la aparición más frecuente.

Sea M la cantidad de documentos en la colección.

Sea k la cantidad de documentos que contienen la palabra w .

Entonces, $IDF(w) = \log[(M + 1)/k]$



IDF (II)

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum_{w \in q \cap d} c(w, q) c(w, d) \log(M + 1/df(w))$$

En vez de usar la cantidad de veces que aparece una palabra en un documento, usar la frecuencia con la que aparece.

Transformación *BM25*:

$$y = (k + 1)x/x + k \text{ con } x = c(w, d)$$

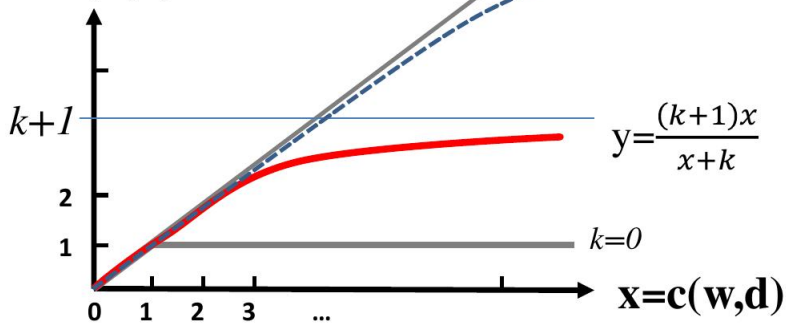
Función de clasificación:

$$f(q, d) = \sum_{i=1}^N x_i y_i$$

$$f(q, d) = \sum_{w \in q \cap d} c(w, q) (k + 1) c(w, d) / c(w, d) + k \log(M + 1) / df(w)$$

Term Frequency Weight

$$y = \text{TF}(w, d)$$



MEJORANDO EPWC

Se pueden utilizar transformaciones sublineales de la frecuencia de términos (TF).

Hay que recompensar a las ocasiones en que coinciden en palabras *raras*.

Se usa IDF.

Se mejora mucho usando BM25.

BM25 Y RELACIONES PARADIGMÁTICAS

Sea $d_1 = (x_1, \dots, x_N)$

$$BM25(w_i, d_1) = (k + 1)(c(w_i, d_1) / (c(w_i, d_1) + k(1 - b + b * |d_1| / avd_1)))$$

Con $b \in [0, 1]$ y $k \in [0, +\infty)$

Entonces $x_i = BM25(w_i, d_1) / \sum_{j=1}^N BM25(w_j, d_1)$

y_i se define de la misma manera.

$$Sim(d_1, d_2) = \sum_{i=1}^N IDF(w_i) x_i y_i$$