

# Motores de Búsqueda

DAI

Un buscador es un programa de *Recuperación de Información*, que trata de buscar las páginas más relevantes a partir de términos de búsqueda

Al contrario que en una base de datos relacional, se busca en **información no estructurada** (texto)

- Una aplicación muy usual en cualquier sitio web
- **Rentable** <http://www.gurusblog.com/archives/google-presenta-resultados-y-subes-casi-un-5-en-el-after-hours/23/01/2013/>
- **Importante, SEO**  
[http://es.wikipedia.org/wiki/Posicionamiento\\_en\\_buscadores](http://es.wikipedia.org/wiki/Posicionamiento_en_buscadores)
- ¿Como sería internet sin buscadores?



# El espacio vectorial

Se asigna a cada página un vector en la que cada coordenada representa la importancia de un posible término de búsqueda (palabra) en la página

El procedimiento sería:

- 1 Se saca un vocabulario con todas las palabras que aparezcan en cualquier página
- 2 Se ordenan por orden alfabético
- 3 A cada palabra se le asigna una dimensión en el espacio vectorial
- 4 A cada página se le asigna un vector, en el que cada coordenada refleje la importancia de la palabra en la página

```
pagina_1 = '<b>Hola</b>, que haces'
```

```
pagina_2 = 'Tienda on-line'
```

```
Vocabulario=['haces', 'hola', 'on-line', 'que',  
             'tienda']
```

```
pag_1 = [1,2,0,1,0] # 'Hola' está resaltada
```

```
pag_2 = [0,0,1,0,1]
```

	<i>haces</i>	<i>hola</i>	<i>on – line</i>	<i>que</i>	<i>tienda</i>
<i>pagina_1</i>	1	2	0	1	0
<i>pagina_2</i>	0	0	1	0	1

La tabla se puede leer en horizontal o en vertical

# Stop words (palabras vacías)

Podríamos quitar palabras sin significado, que no se van a usar en las búsquedas

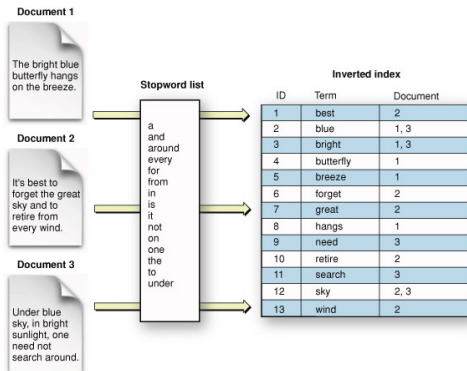
	<i>haces</i>	<i>hola</i>	<i>on – line</i>	<i>tienda</i>
<i>pagina_1</i>	1	2	0	0
<i>pagina_2</i>	0	0	1	1

Así los índices son más pequeños



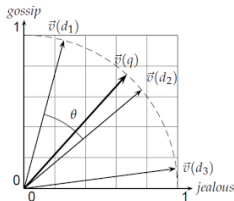
# índice inverso

Lo que hacemos es '*invertir*' el texto; en lugar de archivos de páginas con palabras, tenemos índices de palabras con páginas



Para computar el parecido de un vector a otro usamos el **coseno** del ángulo que forman

$$\text{similitud}(\vec{a}, \vec{b}) = \cos(\theta) = \frac{\sum a_i b_i}{\|\vec{a}\| \|\vec{b}\|}$$



Cosine similarity illustrated.  $\text{sim}(d_1, d_2) = \cos \theta$ .

De esta manera ordenamos los resultados de la búsqueda por relevancia

Consulta: 'Hola tienda' =  $(0, 1, 0, 1)$

Consulta . pagina\_1 =  $(0, 1, 0, 1) \cdot (1, 2, 0, 0) = 2$

Consulta . pagina\_2 =  $(0, 1, 0, 1) \cdot (0, 0, 1, 1) = 1$

# Peso de las palabras en al página

Para el peso de la palabra en la página podemos tener en cuenta:

- El número de veces que se repite la palabra en la página
- El lugar que ocupe (las palabras en los títulos o encabezados serán más importantes)
- El resaltado (palabras en negrita, etc)

En todo caso este número tiene que estar **normalizado**, es decir cada página tiene el mismo peso que reparte entre las palabras que aparecen en ella.

# Peso de las palabras: idf

Inverse Document Frequency

No todas las palabras en una consulta deben tener la misma importancia en la búsqueda.

pe. en la búsqueda 'Restaurantes de Granada'

buscamos páginas en las que sea importante la palabra '**Restaurantes**' y la palabra '**Granada**',  
pero no la palabra 'de'

# Peso de las palabras: idf

Inverse Document frequency

Esto se corrige aplicando un peso a las palabras de la búsqueda **Frecuencia inversa por documento**

$$idf = \log \frac{N}{n_i}$$

$N$  es el número total de páginas

$n_i$  es el número de páginas donde aparece la palabra  $i$

Las palabras comunes, (*stop words*), tendrán un `idf`, cercano a 0, y por tanto podemos quitarlas del índice para ahorrar espacio y tiempo de proceso

Al final tenemos el esquema de pesos tf-idf

$$\text{TF-IDF}_{(n,d)} = \text{TF}_{(n,d)} \times \text{IDF}_{(n)}$$

The diagram illustrates the TF-IDF formula. The formula is  $\text{TF-IDF}_{(n,d)} = \text{TF}_{(n,d)} \times \text{IDF}_{(n)}$ . Below the formula, three components are highlighted with colored brackets and boxes:

- TF-IDF<sub>(n,d)</sub>** (red bracket and box): **Peso de un término (n) en un documento (d)**
- TF<sub>(n,d)</sub>** (blue bracket and box): **Frecuencia de aparición de un término (n) en un documento (d)**
- IDF<sub>(n)</sub>** (yellow bracket and box): **Factor IDF de un término (n)**

<http://en.wikipedia.org/wiki/Tf%E2%80%93idf>

# Seudocódigo

```
relevancia = {}

for t in palabras_busqueda:
    idf = IDF_DE(t)
    lista_pag = PAGS_DE(t)

    for pag in lista_pag:
        relevancia[pag] += PUNT_EN(t, pag) * idf;
    }
```



Otro truco para ahorrar espacio y tiempo en las búsquedas  
indexar el morfema raíz de las palabras, quitando el morfema  
final

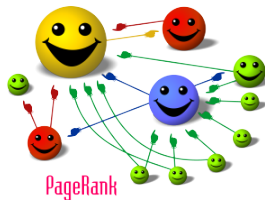
P.e. podemos indexar `niñ`, cuando aparezca *niño*, *niños*, *niña*,  
*o niñas*

Así el índice es más pequeño y al buscar *niño*, también  
recuperamos páginas en las que aparece *niña o niños*

# Peso de las páginas

## Page Rank

También podemos tener en cuenta para el que unas páginas sean más relevantes que otras



Google utiliza un algoritmo en el que se tienen en cuenta los enlaces hacía cada página; una pagina será más importante si otras páginas importantes la enlazan

<http://es.wikipedia.org/wiki/PageRank>

<http://www.mipagerank.com/>