

Tarea del Tema 5: Estudio de diferentes representaciones vectoriales

A partir de un pequeño corpus de páginas web (puedes descargar un conjunto de documentos HTML a tu elección), el objetivo de esta tarea es crear las representaciones de los documentos de dicho corpus dentro del *modelo de espacio vectorial* y usando, para cada uno de ellos, al menos dos de las funciones de pesado presentadas en el tema, una de carácter local y otra global. Deberán entregarse también documentos de texto con la información necesaria para generar las representaciones (el vocabulario, el fichero invertido y las propias representaciones).

Para generar las representaciones deberás aplicar las siguientes fases:

- **Análisis léxico**; seleccionando los rasgos con los que generar las representaciones. Deberás considerar los espacios en blanco como separadores entre cadenas. El resto de decisiones (eliminación de signos de puntuación, caracteres raros, ...) se dejan a tu elección.
- **Eliminación de stop-words**. Deberás buscar una lista de stop-words correspondiente a la lengua en la que estén escritos los documentos que hayas seleccionado y aplicarla a la fase de eliminación de palabras vacías.
- **Truncado**. Aplicar las reglas del algoritmo de Porter para textos en inglés (<http://tartarus.org/martin/PorterStemmer/>). En caso de que prefieras emplear documentos escritos en otra lengua, deberás buscar otro algoritmo de *stemming* y aplicar sus reglas de truncado (<http://snowball.tartarus.org/texts/stemmersoverview.html>).
- **Aplicación de dos funciones de pesado**, una de carácter local y otra global.

Finalmente deberás entregar una pequeña memoria que contenga una sección para cada uno de los puntos previos, y donde se expliquen razonadamente, y en cada caso, las decisiones que has ido tomando.

Nota:

Tanto el Vocabulario como las diferentes representaciones deberán generarse en un documento TXT, donde el formato de representación podrá ser:

- Vocabulario: **term₁ term₂ ... term_N** (en una sola línea y separados por espacios/tab)
- Representaciones: **pesoTerm₁ pesoTerm₂ pesoTerm_N** (en una sola línea y separados por espacios/tab)
- Fichero invertido: **term₁ -> doc₁:TF_{1,1} doc₂:TF_{1,2} doc₃:TF_{1,3} ... doc₁₀:TF_{1,10}**
term₂ -> doc₁:TF_{2,1} doc₂:TF_{2,2} doc₃:TF_{2,3} ... doc₁₀:TF_{2,10}
...

La información por término en una línea y donde **TF_{2,1}** representa la frecuencia de **term₂** en **doc₁**

En caso de usarse otro formato, deberá detallarse claramente en la memoria.