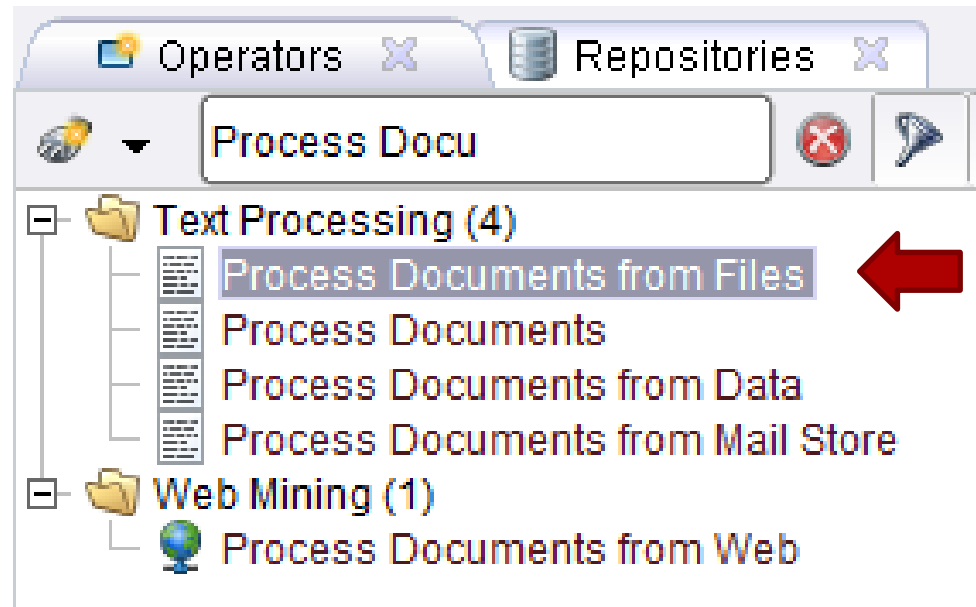




CLUSTERING DE DOCUMENTOS

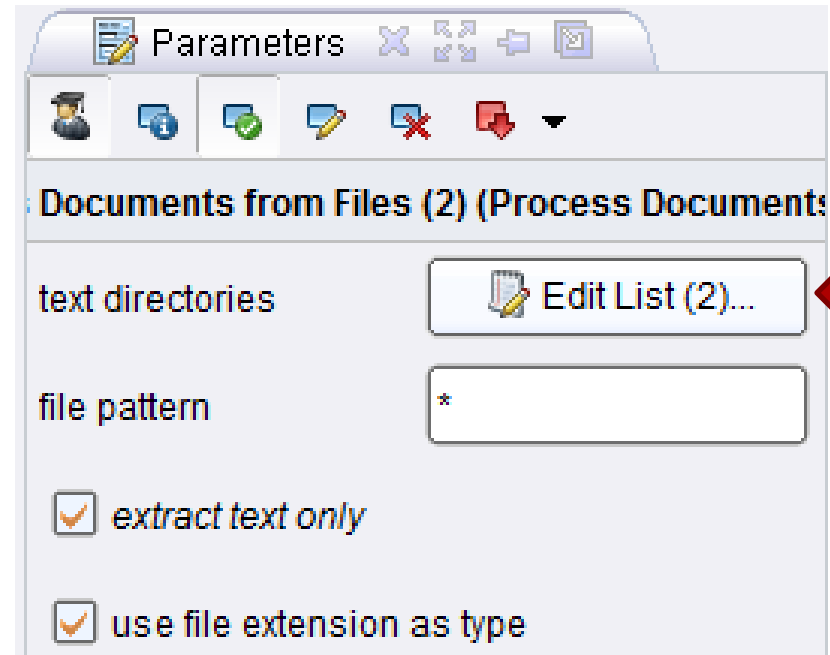
CLUSTERING DE DOCUMENTOS

- Para aplicar una técnica de clustering basada en centroides es preciso tener una representación vectorial numérica de los datos de entrada.



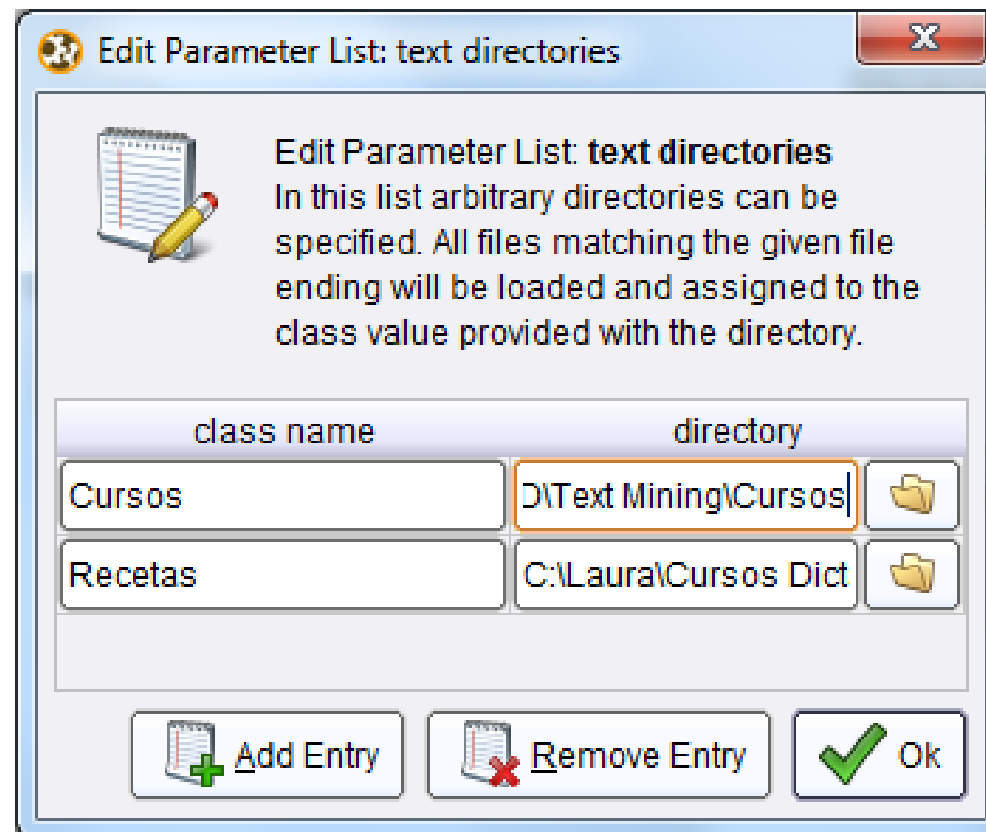
OPERADOR PROCESS DOCUMENTS FROM FILE

- Comencemos indicando la ubicación de los documentos



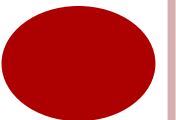
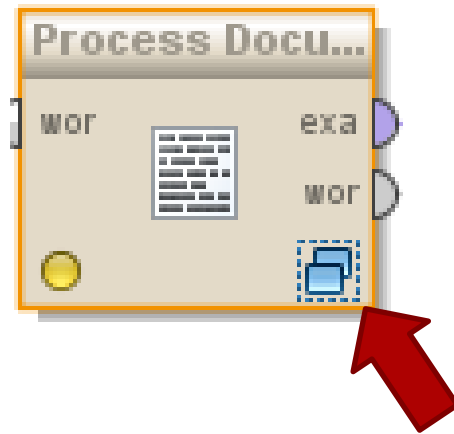
OPERADOR PROCESS DOCUMENTS FROM FILE

- Ubicación de los documentos



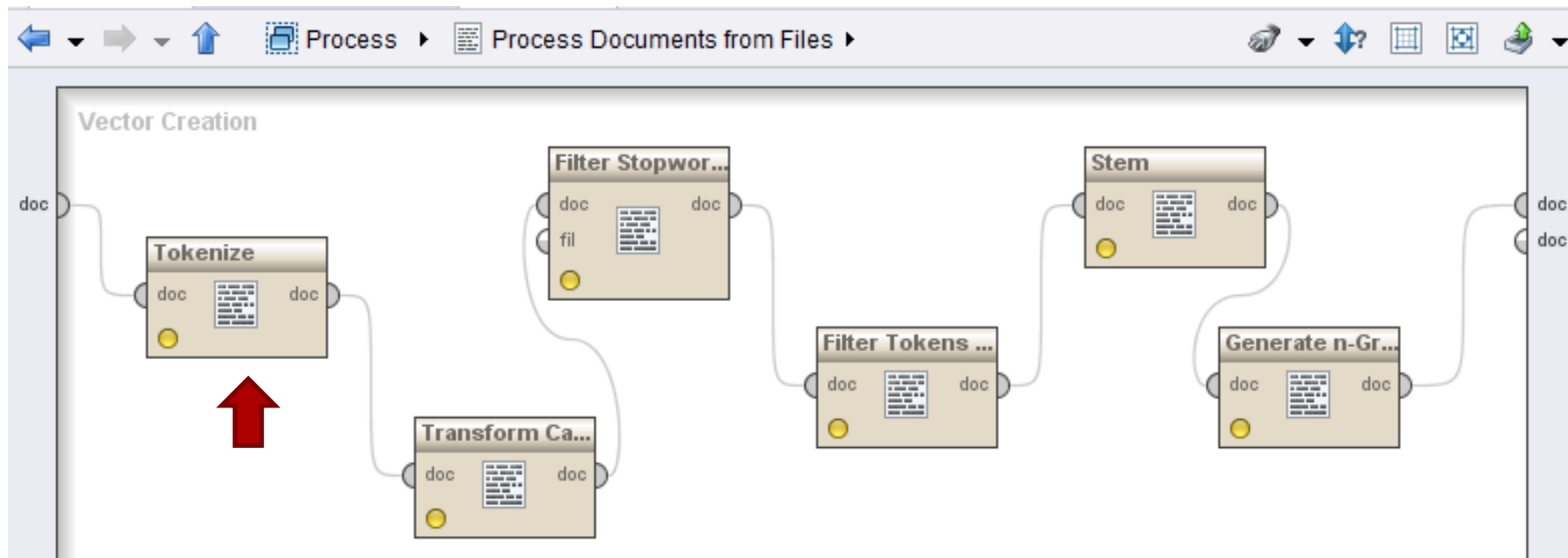
OPERADOR PROCESS DOCUMENTS FROM FILE

- Dentro de este operador se indica la manera de separar cada texto en palabras



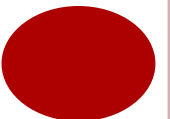
CÓMO SE OBTIENEN LOS TÉRMINOS DE LOS DOCUMENTOS

- Dentro del operador Process Documents from file



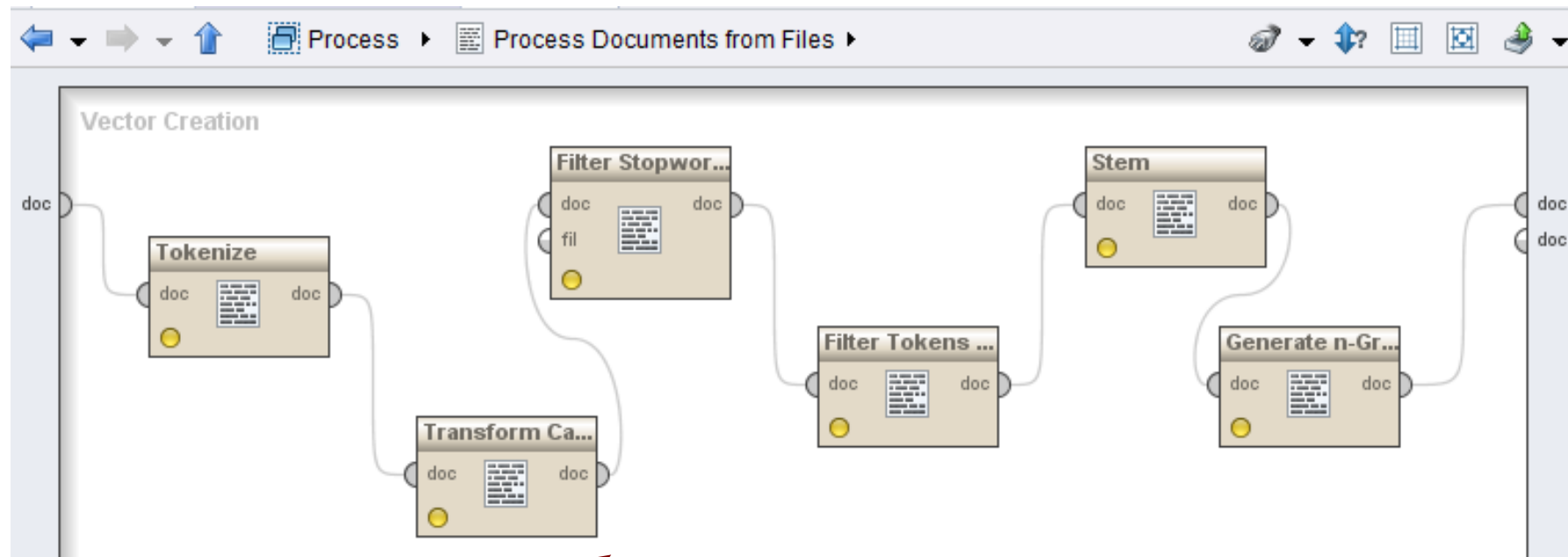
TOKENIZE

Separa cada documento en tokens
(puede indicarse el carácter a utilizar)



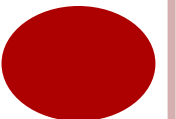
CÓMO SE OBTIENEN LOS TÉRMINOS DE LOS DOCUMENTOS

- Dentro del operador Process Documents from file



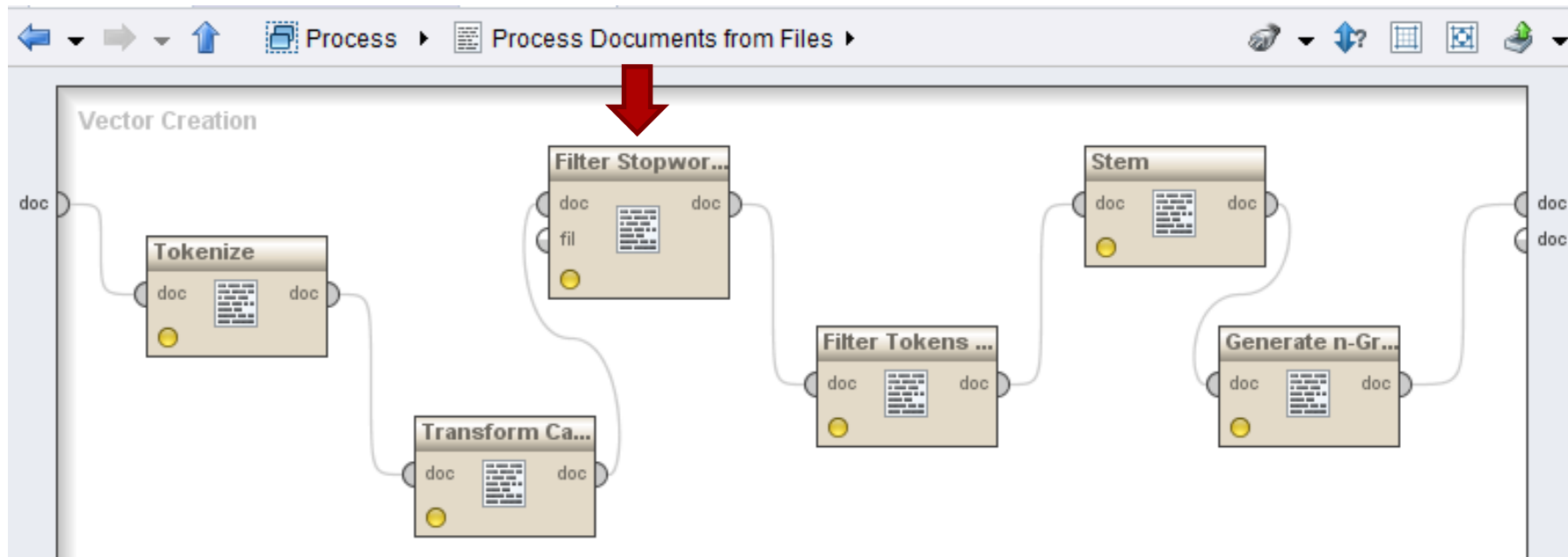
TRANSFORM CASES

Transforma todos los caracteres a minúsculas o mayúsculas según se indique (**seleccionar minúsculas**)



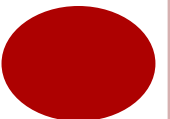
CÓMO SE OBTIENEN LOS TÉRMINOS DE LOS DOCUMENTOS

- Dentro del operador Process Documents from file



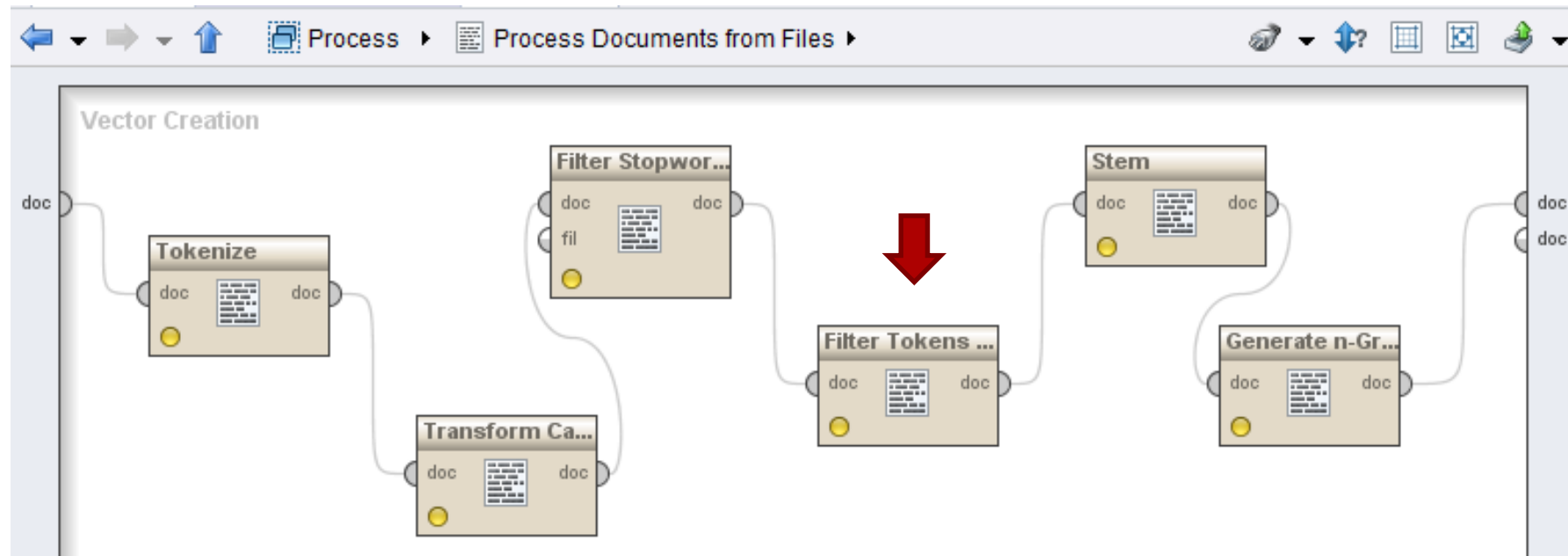
FILTER STOPWORDS (DICTIONARY)

Filtra los tokens que coincidan con cualquier stopwords indicada en un determinado archivo (utilice el archivo **stopwords_es.txt**)

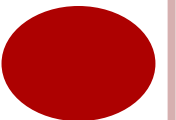


CÓMO SE OBTIENEN LOS TÉRMINOS DE LOS DOCUMENTOS

- Dentro del operador Process Documents from file

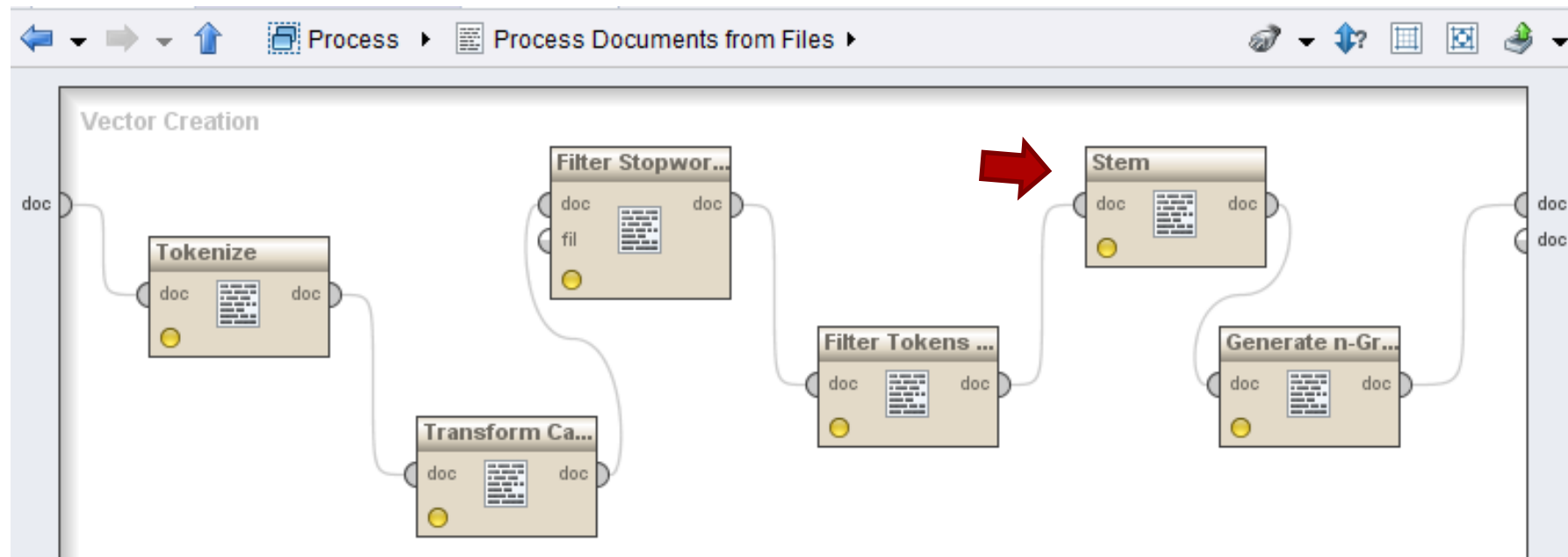


FILTER TOKENS (by LENGTH)
Configúrelo para que utilice sólo los tokens
entre 4 (min chars) y 25 (max chars) caracteres.



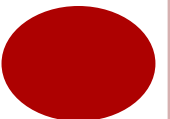
CÓMO SE OBTIENEN LOS TÉRMINOS DE LOS DOCUMENTOS

- Dentro del operador Process Documents from file



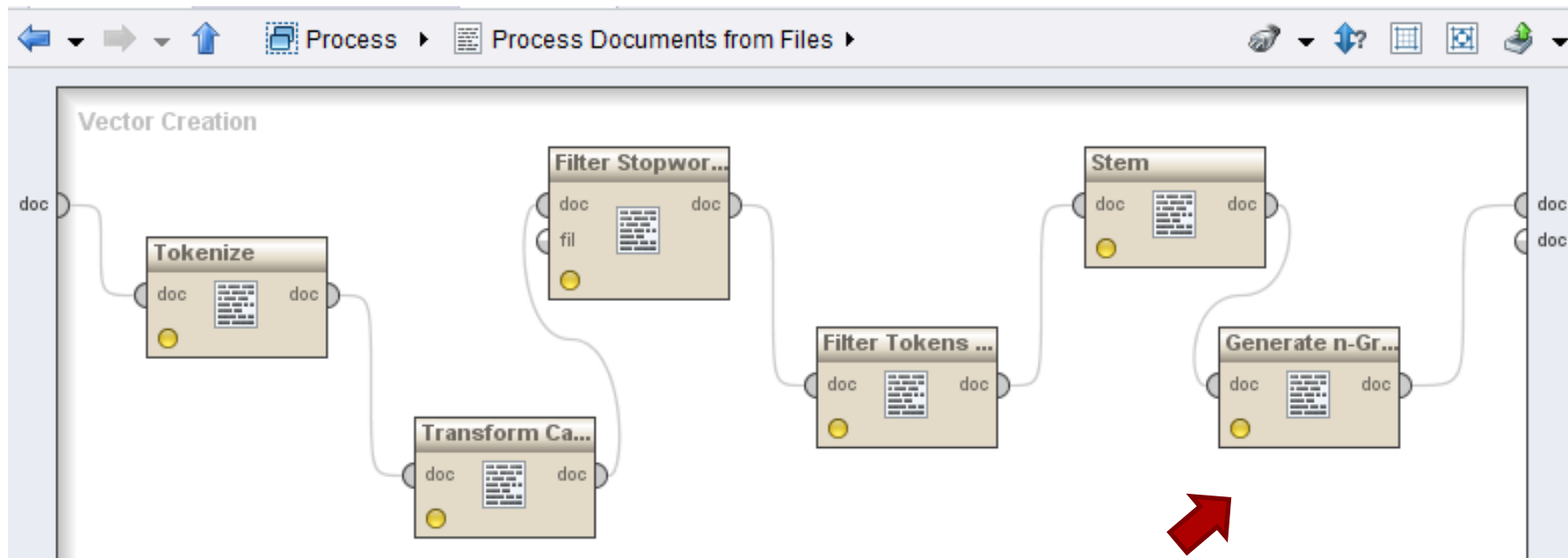
STEM (SNOWBALL)

Aplica un algoritmo de stemming para el lenguaje seleccionado.

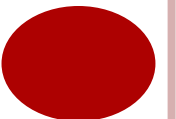


CÓMO SE OBTIENEN LOS TÉRMINOS DE LOS DOCUMENTOS

- Dentro del operador Process Documents from file

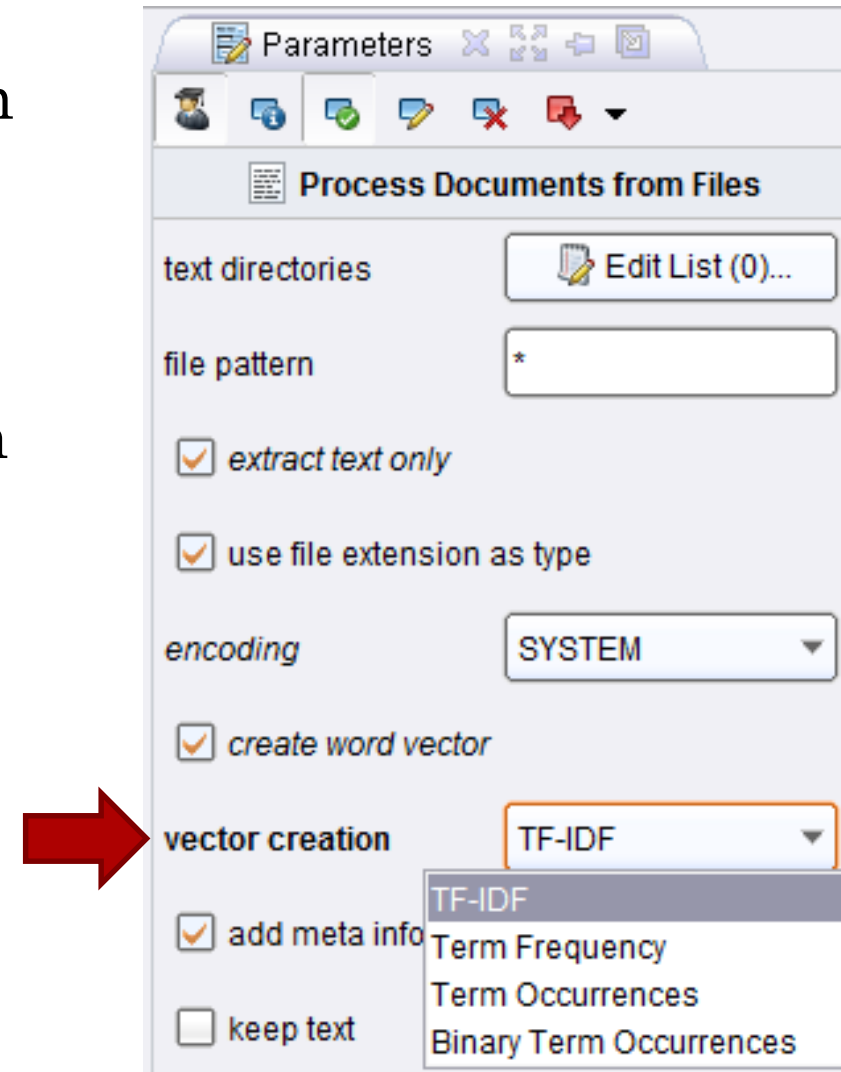


Grenerate n-Grams (terms)
Combina tokens. Utilice max length = 2



REPRESENTACIÓN NUMERICA

- Se generará un gran diccionario con las palabras de todos los documentos.
- Es preciso indicar la representación numérica a utilizar



FRECUENCIA DE UN TÉRMINO

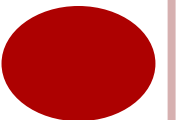
Opciones

- Si es binaria, $tf(t, d) = 1$ si el término t aparece en el documento d y 0 en caso contrario.
- Si la escala es logarítmica

$$tf(t, d) = \begin{cases} 1 + \log(f(t, d)) & f(t, d) > 0 \\ 0 & f(t, d) = 0 \end{cases}$$

- La frecuencia de un término puede normalizarse utilizando

$$tf(f, d) = \frac{f(t, d)}{\max\{f(t, d): t \in d\}}$$



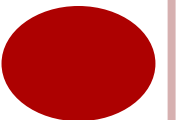
INVERSE DOCUMENT FREQUENCY

- Es una medida de cuán común o raro es un término entre todos los documentos

$$idf(t, D) = \log \left(\frac{|D|}{|\{d \in D : t \in d\}|} \right)$$

siendo

- $|D|$ el número total de documentos.
- $|\{d \in D : t \in d\}|$ el número de documentos en los que aparece el término t . Puede evitar la división por 0 usando $1 + |\{d \in D : t \in d\}|$
- La base del logaritmo es sólo un factor constante multiplicativo.

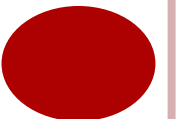


TF-IDF

- Se calcula de la siguiente forma

$$tf - idf(t, d, D) = tf(t, d) * idf(d, D)$$

- Se obtendrán valores altos en aquellos términos que posean frecuencia alta (en el documento dado) con una frecuencia baja del término en la colección completa de documentos, filtrando de esta forma los términos comunes.
- Dado que la proporción usada como argumento de la función *log* en *idf* es > 0 , *tf - idf* también es > 0 siempre.
- Un término que aparece en muchos documentos hace que la proporción del logaritmo en *idf* se acerque a 1 y por lo tanto *idf* y *tf - idf* tenderán a 0.



AGRUPAMIENTO

