

Assignment No. 4

Title: Consider a suitable text dataset. Remove stop words, apply stemming & feature & feature selection techniques to represent documents as vectors. Classify documents & elevate precision, recall.

Problem Definition:-

Remove stop words.

Pre-requisite:-

Basic concepts of ETL.

S/W & H/W Requirements:-

Rapid Miner, PIV, 2GB RAM, 500 GB HDD, Lenovo A13-4089 Model.

Learning Objectives:-

We are going to learn how to tokenize & filter a document into its different words & then do words count for each word in a text document.

Outcomes:-

You are able to see a word list containing all the different words in your document & their occurrence count next to it in the

"Total Occurrences" column.

Theory Concepts:-

Text Processing Tutorial with Rapid Miner:

In this manual, we are going to learn how to tokenize & filter a document into its differentiate words & then do words count for each word in a text document.

Open RapidMiner & click 'New Process'. On the left hand panel of your screen, there should be a tab that says 'Operators' - this is where you can search & find all the operators for Rapid Miner & its extensions. By searching the Operators tab for 'read', you should get an output like this. There are multiple read operators depending on which file you have & most of them work the same way. If you have & most of them work the same way. If you scroll down, there is a 'Read Document' operator. Select this operator & enter it into your Main Process window by dragging it. When you select the Read Documents operator in the Main Process Window, you should see a file uploader in the right.

After you have chosen your file, make sure that the output port on the Read Documents operator is connected to the 'res' node in your Main process. Click

the 'play' button to check that your Main process. Click the 'play' button to check that your file has been received correctly. Switch to the results perspective by clicking the icon that looks like a display chart above the 'Processes' tab at the top of the Main Process panel.

Now we will move on processing the ~~document~~ document to get a list of its different words & their individual count. Search the Operator list for 'Process Document' to get a list of its different words & their individual count. Search the Operator list for 'Process Documents'.

Double Click the Process Documents operator to get inside the operator. This is where we will link operators together to take the entire text document & split it down into its word components. This consists of several operators that can be chosen by going into the Operator panel & ~~are~~ looking at the Text Processing folder.

Now we are ready to filter the bag of words. In 'filtering' folder under the 'Text Processing' operator folder, you can see the various filtering methods that you can apply to your process.

Conclusion:-

We are now able to see a word list containing all the different words in your document & their occurrence count next to it in the 'Total Occurrence' column. If you do not get this output, make sure that all of your nodes are connected correctly & also to the right type. Some errors are because your output at one node does not match the type expected at the i/p of the next node of an operator.



Views:

Design

Results

Turbo Prep

Auto Model

Deployments

Find data, operators...etc



All Studio

Repository

Import Data

- Samples
- Local Repository (Local)
 - Connections
 - data
 - processes
 - Ass1 | 10/13/21 2:47 PM - 1 KB
 - ass4_text | 10/17/21 8:39 PM -
 - kmean_ass2 | 10/15/21 6:22

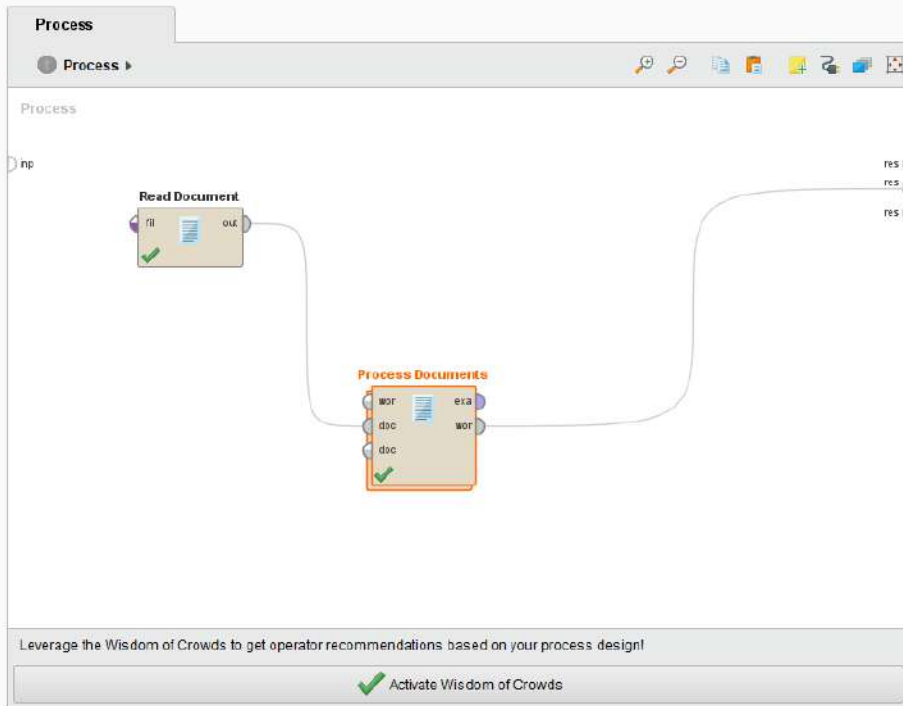
Operators

transform cases

Extensions (1)

- Text Processing (1)
 - Transformation (1)
 - Transform Cases

No results were found.



Parameters

Process Documents

- ☒ create word vector
- vector creation: TF-IDF
- ☒ add meta information
- ☐ keep text
- prune method: none
- data management: auto

[Hide advanced parameters](#)

Help

Process Documents

Text Processing

Tags: Text Processing

Synopsis

Generates word vectors from a text object.

Description



Views:

Design

Results

Turbo Prep

Auto Model

Deployments

Find data, operators...etc

All Studio

Repository

Import Data

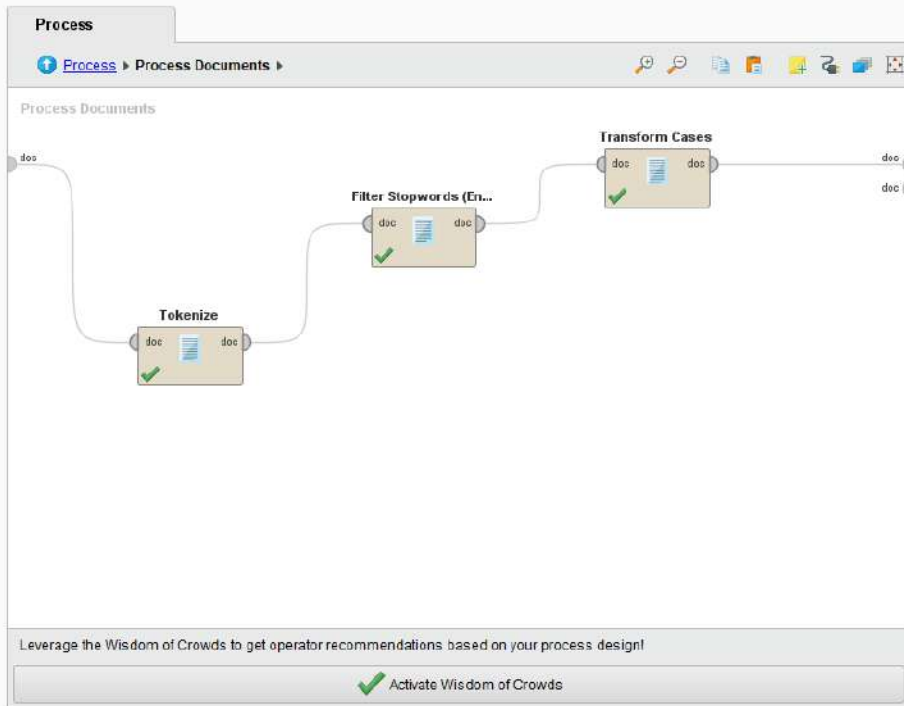
- Samples
- Local Repository (Local)
 - Connections
 - data
 - processes
 - Ass1 | 10/13/21 2:47 PM - 1 KB
 - ass4_text | 10/17/21 8:39 PM -
 - kmean_ass2 | 10/15/21 6:22

Operators

transform cases

- Extensions (1)
 - Text Processing (1)
 - Transformation (1)
 - Transform Cases

No results were found.



Parameters

Process Documents

- ☒ create word vector
- vector creation: TF-IDF
- ☒ add meta information
- ☐ keep text
- prune method: none
- data management: auto
- [Hide advanced parameters](#)

Help

Process Documents

Text Processing

Tags: Text Processing

Synopsis

Generates word vectors from a text object.

Description