

# tech\_review

November 15, 2020

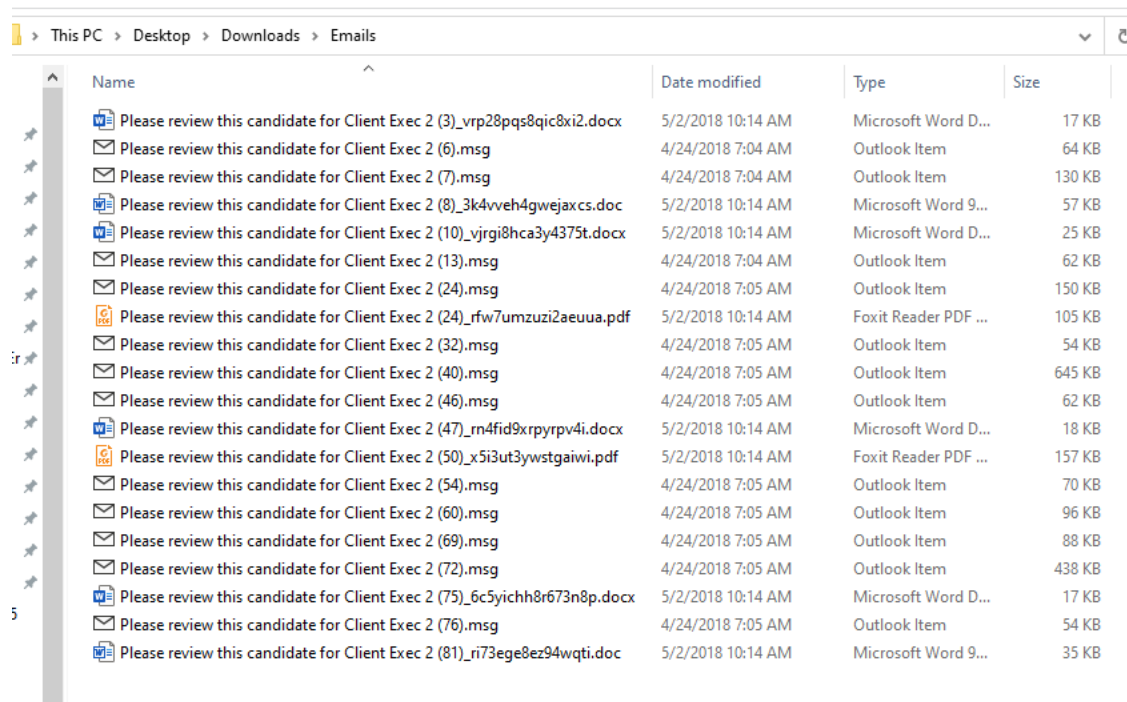
## 1 TF/IDF with KNIME

In this tech review, I am going to show how to compute TF/IDF with KNIME. KNIME is a drag-n-drop advanced analytics platform. More at [KNIME.com](https://www.knime.com)

KNIME Desktop is a free software. It can do all of the analytics that one might need to do.

In this demo, I am going to show how to parse documents (resumes) and compute TF/IDF.

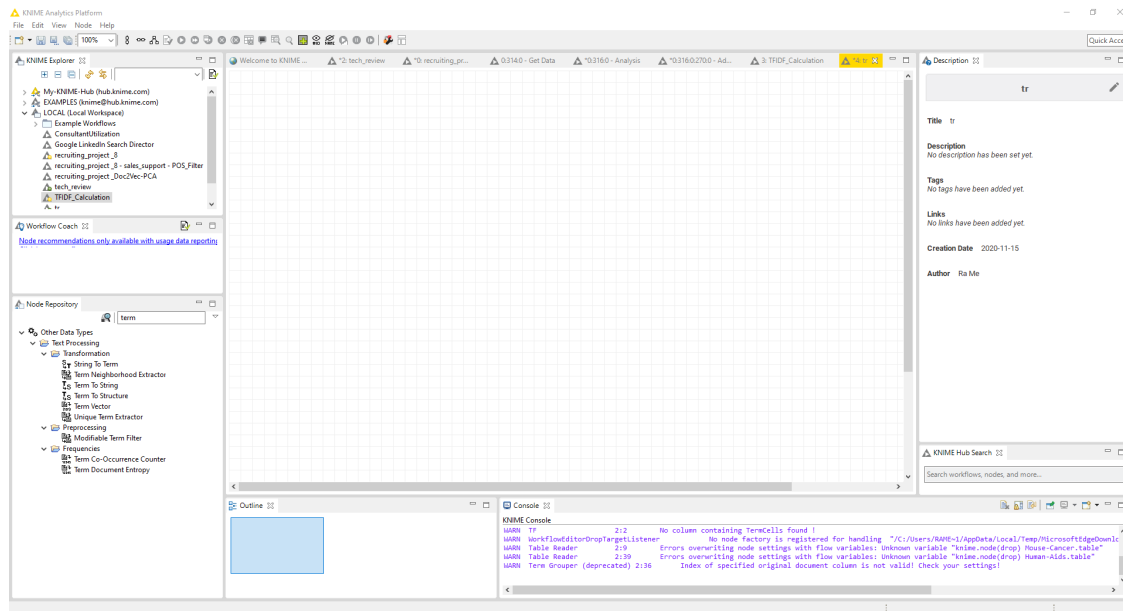
I have a folder with resumes in it that looks like this:



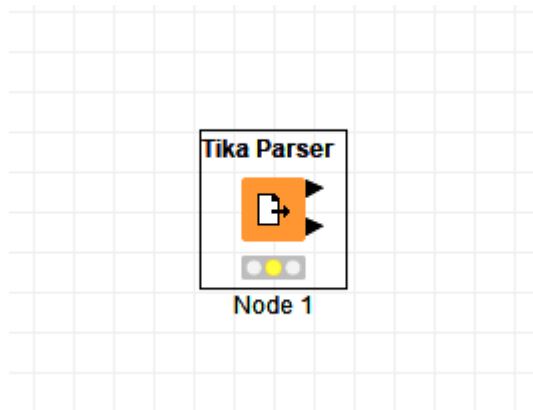
Name	Date modified	Type	Size
Please review this candidate for Client Exec 2 (3)_vrp28pqs8qic8xi2.docx	5/2/2018 10:14 AM	Microsoft Word D...	17 KB
Please review this candidate for Client Exec 2 (6).msg	4/24/2018 7:04 AM	Outlook Item	64 KB
Please review this candidate for Client Exec 2 (7).msg	4/24/2018 7:04 AM	Outlook Item	130 KB
Please review this candidate for Client Exec 2 (8)_3k4vveh4gwejajcs.doc	5/2/2018 10:14 AM	Microsoft Word 9...	57 KB
Please review this candidate for Client Exec 2 (10)_vjrgi8hca3y4375t.docx	5/2/2018 10:14 AM	Microsoft Word D...	25 KB
Please review this candidate for Client Exec 2 (13).msg	4/24/2018 7:04 AM	Outlook Item	62 KB
Please review this candidate for Client Exec 2 (24).msg	4/24/2018 7:05 AM	Outlook Item	150 KB
Please review this candidate for Client Exec 2 (24)_rfw7umzuzi2aeuua.pdf	5/2/2018 10:14 AM	Foxit Reader PDF ...	105 KB
Please review this candidate for Client Exec 2 (32).msg	4/24/2018 7:05 AM	Outlook Item	54 KB
Please review this candidate for Client Exec 2 (40).msg	4/24/2018 7:05 AM	Outlook Item	645 KB
Please review this candidate for Client Exec 2 (46).msg	4/24/2018 7:05 AM	Outlook Item	62 KB
Please review this candidate for Client Exec 2 (47)_m4fid9xrpyrpv4i.docx	5/2/2018 10:14 AM	Microsoft Word D...	18 KB
Please review this candidate for Client Exec 2 (50)_x5i3ut3ywtgaiwi.pdf	5/2/2018 10:14 AM	Foxit Reader PDF ...	157 KB
Please review this candidate for Client Exec 2 (54).msg	4/24/2018 7:05 AM	Outlook Item	70 KB
Please review this candidate for Client Exec 2 (60).msg	4/24/2018 7:05 AM	Outlook Item	96 KB
Please review this candidate for Client Exec 2 (69).msg	4/24/2018 7:05 AM	Outlook Item	88 KB
Please review this candidate for Client Exec 2 (72).msg	4/24/2018 7:05 AM	Outlook Item	438 KB
Please review this candidate for Client Exec 2 (75)_6c5yichh8r673n8p.docx	5/2/2018 10:14 AM	Microsoft Word D...	17 KB
Please review this candidate for Client Exec 2 (76).msg	4/24/2018 7:05 AM	Outlook Item	54 KB
Please review this candidate for Client Exec 2 (81)_ri73ege8ez94wqti.doc	5/2/2018 10:14 AM	Microsoft Word 9...	35 KB

You can see that we have 20 documents: .MSG, .DOC, .DOCX, and .PDF.

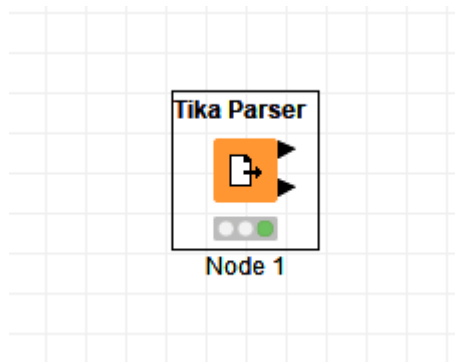
First, we start with a blank workbench.



We want to parse all of our documents at once. In KNIME, there is a processing node called **Tika Parser**. We can drop it into our workbench.



The next step is to specify where the documents are. For that, we right-click, click **Configure**, then click on **Browse**, then navigate to the folder where the documents are and click **OK**. After that, right-click on the node and click **execute**. You will see that the yellow collar under the node will turn to green.



If we'll right-click on the node again, and if we'll click on **Metadata output table**, we'll see a table with parsed data. The last column, **Content**, contains the actual document data.

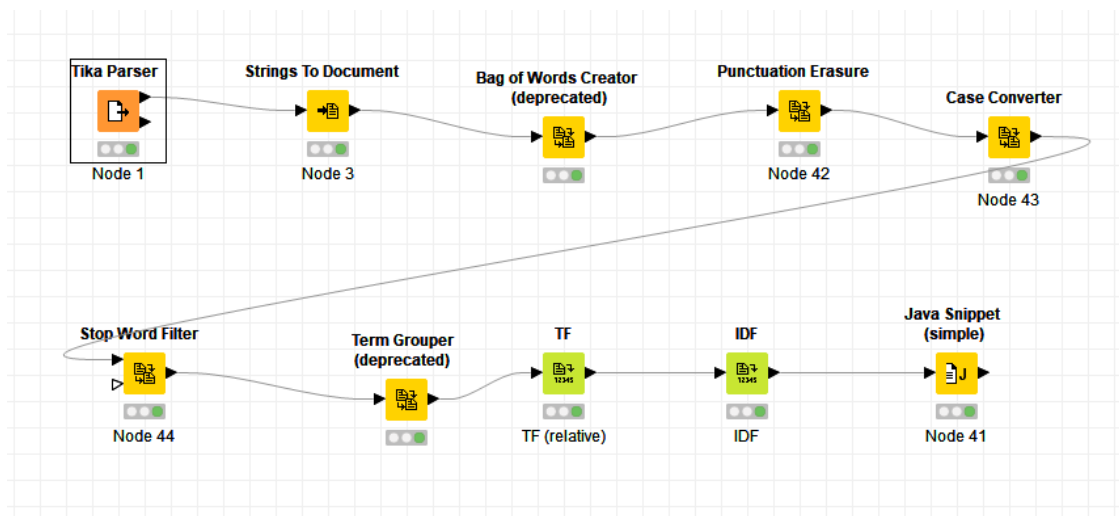
Metadata output table - #1 - Tika Parser  
File | Hints | Navigation | View

Table 'default' - Rows: 20 | Spec: Columns: 25 | Properties | Flow Variables

Row ID	Version	Creator	Comment	Metadata	Content
Row0					JOSEPH FRELINO 1155 Oriole Lane Saint Louis, Missouri 63112 (314) 206-1291 frelino1552@yahoo.com
Row1					Please review this candidate for: Client Exec 2 From: jbrown@browner.com To: Walling, Melissa (STL) Recipients:
Row2					Please review this candidate for: Client Exec 2 From: jbrown@browner.com To: Walling, Melissa (STL) Recipients:
Row3		Microsoft			Justin M. Brown www.linkedin.com/in/justin-brown-404ba66 brown.justin.mitchell@gmail.com • 307-369-4152  RELEVANT WORK EXPERIENCE  Material Control Systems Analyst March 2017- Present Duke Manufacturing St. Louis, Missouri • Reduced total inventory value on hand by \$250,000 and reduced stock-out issues by implementing and managing a Kanban Just-in-Time inventory reordering system on four major production lines as a plant leader in Lean Manufacturing Continuous Improvement. • Negotiated lower Minimum Order Quantities and shorter Lead Times with suppliers, and calculated and implemented proper safety stocks and reorder points in order to lower inventory value kept on hand. • Managed new part setup and existing part maintenance, inventory control and warehouse setup, buyer setup and permissions, supplier setup, and supplier price lists in SAP/ ERP system.  Inventory Buyer/Production Scheduler January 2017- March 2017 Duke Manufacturing St. Louis, Missouri • Analyzed inventory levels, scheduled blanket orders, and purchased over \$50,000 in spend per month for parts, packaging, and supplies on 3 major production lines, leading to zero stock-outs and zero line shut-downs. • Communicated daily with suppliers both domestic and international regarding pricing, lead times, and past-due shipments.  Deborah A. Isison Phone: 314-793-2773 daisison@gmail.com www.linkedin.com/in/deborah-isison-3056175 I have extensive experience working for a Fortune 500 company. My career objective is to utilize my experience with the job that match skills. My strong work ethic includes, active listening skills, ability to multi-task, solid interpersonal relations, timely completion of tasks, attention to detail, ability to work independ...
Row4					Please review this candidate for: Client Exec 2

You can see that each cell in the **Content** column is a blob of text that contains resume data.

Here is the final workflow so that we can see what we are trying to accomplish:



From the **Tika Parser** node, we are going to the **Strings To Document** node. This is needed because we have to feed a document to the BoW node. This is the way KNIME works. When we right-click on the BoW node and click on **Documents Output Table**, we can see that our row count is now in thousands because for each term in a document we have a row.

Rows: 8920 Spec - Columns: 2 Prop		
	<b>T</b> Term	
	finalist[]	"J
	TRULASKE[]	"J
	COLLEGE[]	"J
	Columbia[]	"J
	1999-2004[]	"J
	Bachelor[]	"J
	Science[]	"J
	Major[]	"J
	:[]	"J
	Banking[]	"J
	Financial[]	"J
	Management[]	"J
	Association[]	"J
	Vice-President[]	"J
	Treasurer[]	"J
	Student[]	"J
	Council[]	"J
	Phi[]	"J
	Kappa[]	"J
	Fraternity[]	"J
	Rush[]	"J
	Chairman[]	"J
	EXPERIENCE[]	"J
	Solid[]	"J
	Gold[]	"J
	Pet[]	"J
	2017-2017[]	"J
	Senior[]	"J
	Analyst[]	"J
	Developed[]	"J
	sales[]	"J
	marketing[]	"J
	analysis[]	"J
	for[]	"J
	our[]	"J
	domestic[]	"J
	international[]	"J

We then use the nodes **Punctuation Erasure**, **Case Converter**, and **Stop Word Filter** to do what we have learn in our lectures when we were working with tokenizers.

The next node is a **Term Grouper** node. This node groups all terms of a document by their text and deletes all tags.

The final 3 nodes are **TF**, **IDF**, and **Java Snippet**. All 3 nodes compute a column for our table. The last node computes a product of  $TF * IDF$ .

When we right-click on the **Java Snippet** node and click on **Appended Table**, we see the following table:

[illegible]

We can see that KNIME has computed what we needed.

This concludes our demo.