

Intro to NLP

→ target

how to PNL

int ↕

Where do we see the application around us

virtual assistant

QnA → IBM Watson

text → google lens

Voice search

Email

translate

Sentiment analysis

information extraction

NLP → Austerio S

sentime

int ↕

Natural Lang

→ T0CEN

↑? Computer → NLU

France → French

? ↕ Transl

human intuition

codify

rules

algo

comp → binary

0/1

Topic model

NLS

google news

DOCUMENT CLUSTERING

- Has applications in knowledge management and information retrieval
- Makes it easy to group similar documents into meaningful groups.

PATTERN IDENTIFICATION

Process of searching large amounts of text for patterns and recognizing features

Helps extract features such as phone numbers and email addresses



PRODUCT INSIGHTS

- Helps to extract large amounts of text, such as customer reviews
- Can reveal insights about products

amazon
flipkart

SECURITY MONITORING

- Monitors and extracts information from news articles and reports for national security purposes

agencies

Text mining allows you to make better informed decisions, automate information-intensive processes, gather business critical insights, and mitigate operational risks.

Review Content

-ve I like it....I can't put it down....but after buying the full set with Rs2110....the 4th book's 3 pages were out and the 5th book which should be the most thick came out to be less thick than the 4th one and the 4th book condition was so bad that I had to put cellotape

The books are fine but definitely not hardcover and what it shows like in the picture. Feels like a scam.

-ve The books are great! However, the font is too small and if put to regular reading, the cover starts wearing off - mine did in a month for Order of Phoenix. Also, as the photo shows, one of the pages of the books is missing - Deathly Hallows, page 17. Make sure you check once you've purchased the books! Lol disappointed.

~ = 100% X

As everyone knows, this is the popular set of books. I bought this as a gift to my daughter. She is very happy about it. No need to review in detail.

+ve I love harry Potter series and this book was actually my turning point in reading not that I did not read before that but I was attracted to series and trilogy after reading this ...now coming to the shipping...it was amazing since it came in one day time although I received my Order of Phoenix with a slight cut but it does look noticeable but never mind

These books came in perfect condition.. But the box had a little tear which is not a problem as the box might have been thrown a bit in the delivery so overall it is a good product but it can be made better 😊

Absolutely loved it ,Good condition book with catchy illustration, Don't read the negative reviews just buy it you will love it 😍

+ve

Packing is terrible. However it's good condition 😊 😊

-ve Till now the service has been excellent....no chance of a single complaint. But not now. I bought this ₹8000+ costing box set of Harry Potter. The outer box was broken....and all the books were cellotaped to stay together. The first two parts of Harry Potter look as if they might have been previously used by someone...there are those marks on the main cover. Secondly, someone has scribbled lines with blue ballpen on the last page of The Chamber of Secrets. It has also got the used marks on the cover. It's completely unexpected. I did not feel like returning it all...as the rest of the books are in good condition. But it does cross my mind that I could have gone for a cheaper set.

It would've been a terrible mistake if I hadn't purchased these books for reading because they were very popular and I thought I would read em. Harry Potter... You don't know anything If you haven't read it

★★★★★ Not hardcover.

Reviewed in India on 7 December 2019

Verified Purchase

The books are fine but definitely not hardcover and what it shows like in the picture. Feels like a scam.



★★★★★ Good and evil:eternal battle.....

Reviewed in India on 7 October 2018

Verified Purchase

I like it....I can't put it down....but after buying the full set with Rs2110....the 4th book's 3 pages were out and the 5th book which should be the most thick came out to be less thick than the 4th one and the 4th book condition was so bad that I had to put cellotape

★★★★★ No need of a review on this...

Reviewed in India on 13 July 2020

Verified Purchase

As everyone knows, this is the popular set of books. I bought this as a gift to my daughter. She is very happy about it. No need to review in detail.



★★★★★ Missing pages and cover wearing off

Reviewed in India on 9 October 2018

Verified Purchase

The books are great! However, the font is too small and if put to regular reading, the cover starts wearing off - mine did in a month for Order of Phoenix. Also, as the photo shows, one of the pages of the books is missing - Deathly Hallows, page 17. Make sure you check once you've purchased the books! Lol disappointed.

★★★★★ Nice books

Reviewed in India on 21 February 2022

Verified Purchase

These books came in perfect condition.. But the box had a little tear which is not a problem as the box might have been thrown a bit in the delivery so overall it is a good product but it can be made better 😊



★★★★★ Loved the series

Reviewed in India on 31 January 2018

Verified Purchase

I love Harry Potter series and this book was actually my turning point in reading not that I did not read before that but I was attracted to series and trilogy after reading this ...now coming to the shipping...it was amazing since it came in one day time although I received my Order of Phoenix with a slight cut but it does look noticeable but never mind



★★★★★ Loved it

Reviewed in India on 26 August 2020

Verified Purchase

Absolutely loved it ,Good condition book with catchy illustration, Don't read the negative reviews just buy it you will love it 😍



★★★★★ Partly satisfied

Reviewed in India on 15 July 2021

Verified Purchase

Till now the service has been excellent....no chance of a single complaint. But not now. I bought this ₹8000+ costing box set of Harry Potter. The outer box was broken....all the books were cellotaped to stay together. The first two parts of Harry Potter look as if they might have been previously used by someone...there are those marks on the main cover. Secondly, someone has scribbled lines with blue ballpen on the last page of Chamber of Secrets. It has also got the used marks on the cover. It's completely unexpected. I did not feel like returning it all...as the rest of the books are in good condition. But it does cross my mind that I could have gone for a cheaper set.

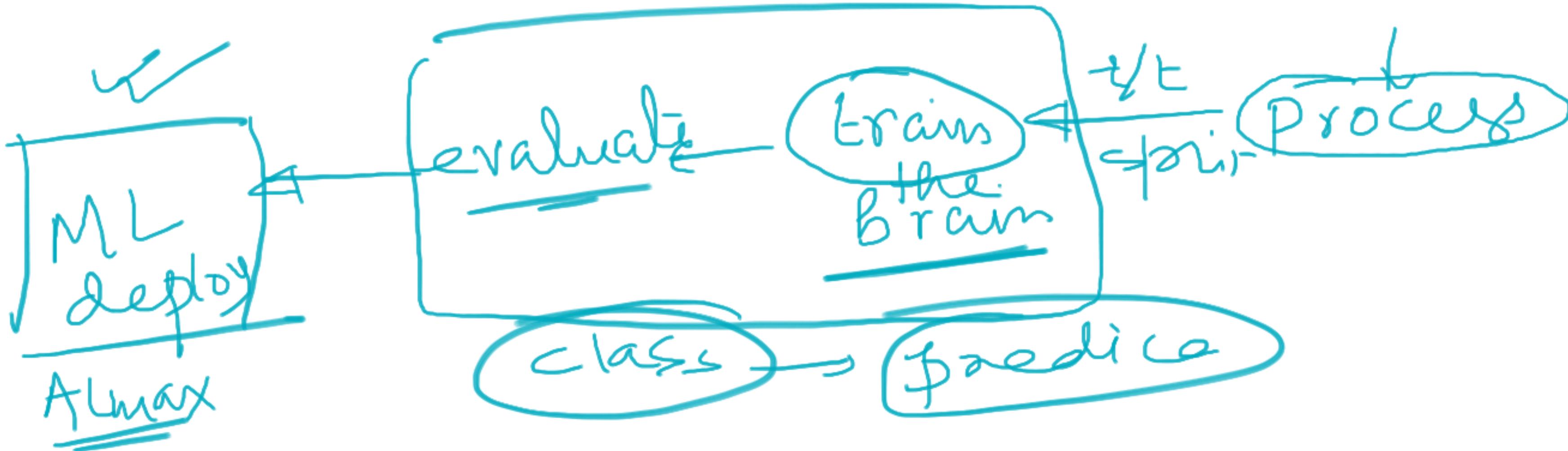
★★★★★ Could tail india terrible packing, cover is half open

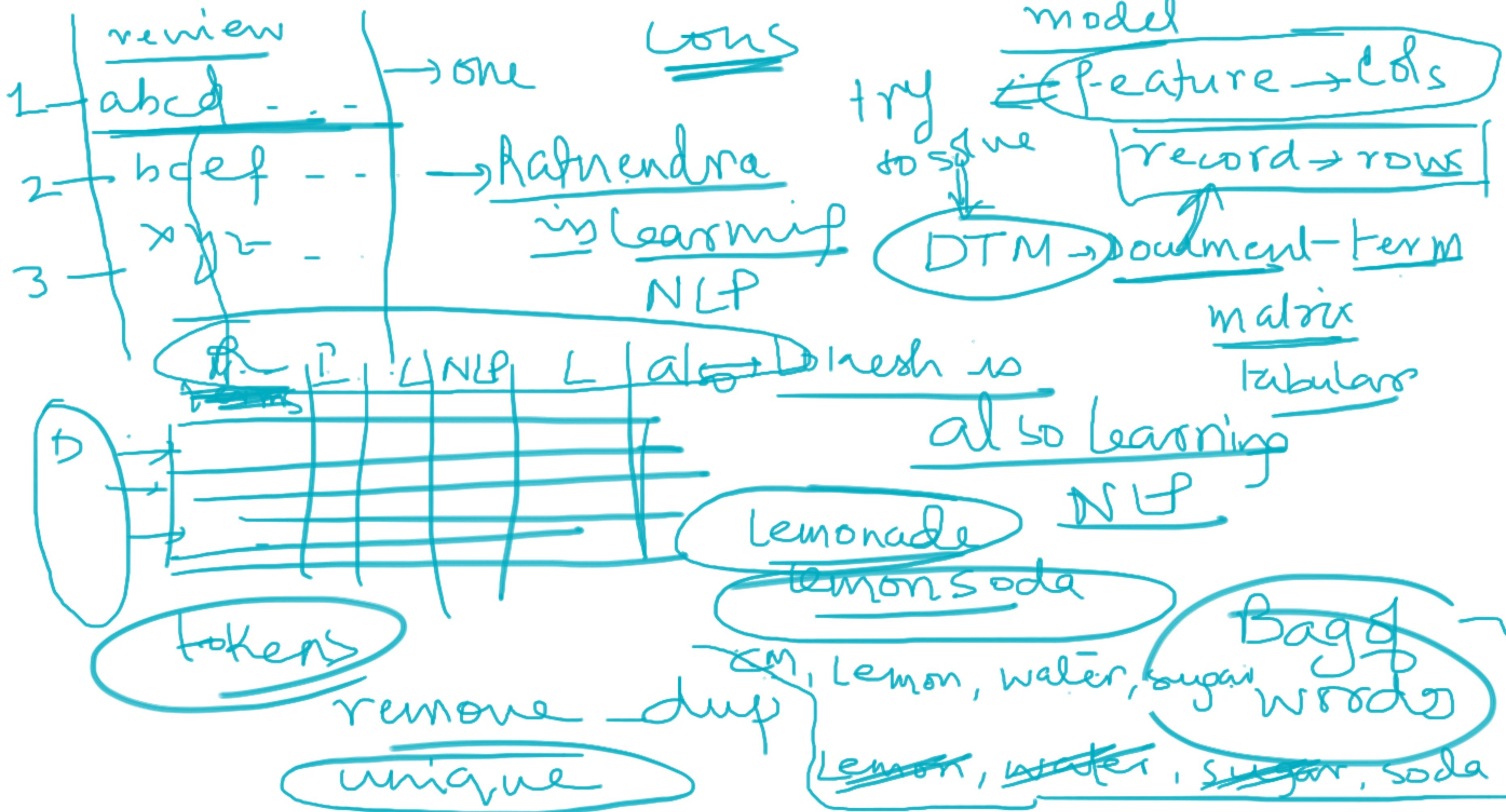
Reviewed in India on 11 September 2021

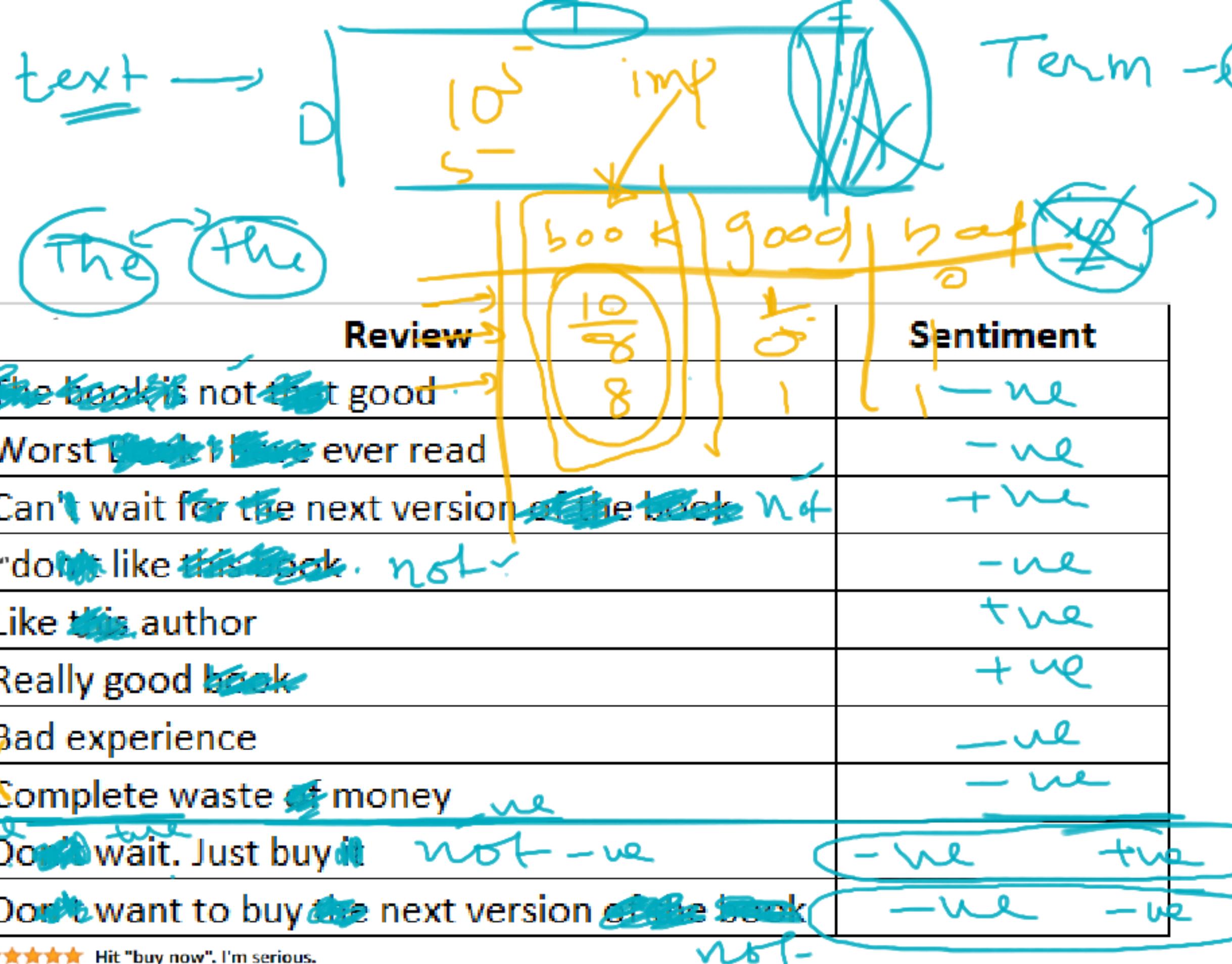
Verified Purchase

Packing is terrible. However it's good condition 😊😊

It would've been a terrible mistake if I hadn't purchased these books for reading because they were very popular and I thought I would read em. Harry Potter... You don't know anything If you haven't read it







 Hit "buy now". I'm serious.

Reviewed in India on 24 December 2016

Verified Purchase

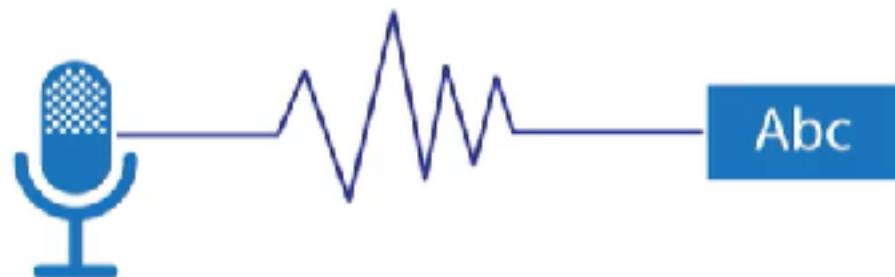
It would've been a terrible mistake if I hadn't purchased these books for reading because they were very popular and I thought I would read em. Harry Potter...

You don't know anything If you haven't read it

synon
+ stemming y Lenne

Applications

SPEECH RECOGNITION



SPAM FILTERING



SENTIMENT ANALYSIS



- It determines if a given sentence expresses positive, neutral, or negative sentiments.
- This includes data analysis of the body of text to understand the opinion, modality, and mood.
- It works best on text that has subjective context.

E-COMMERCE PERSONALIZATION



- Text mining is used to suggest products that fit into a user's profile.
- It helps e-commerce retailers learn more about consumers by analyzing information, identifying patterns, and providing insights.
- Retailers can target specific individuals with personalized offers and discounts.
- It increases customer loyalty by identifying purchase patterns and opinions on products.

~~invent Re-inven~~

~~Stemming~~

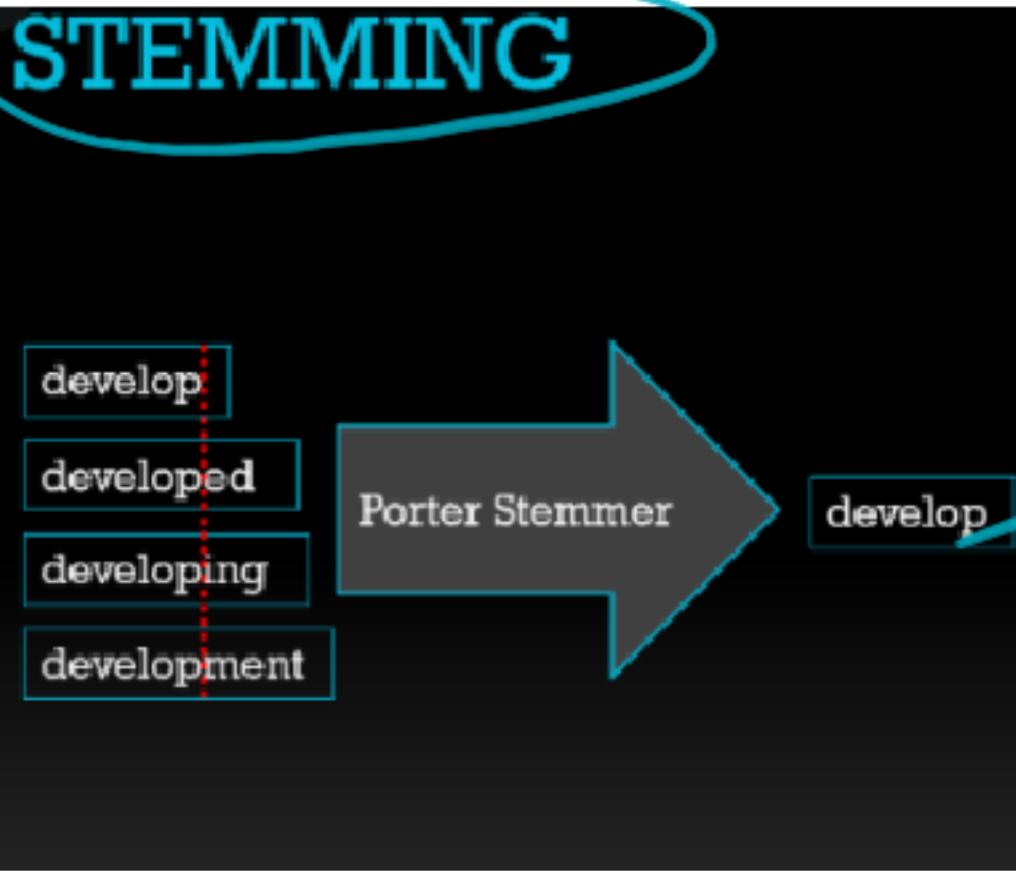
~~Stemming~~

grammati

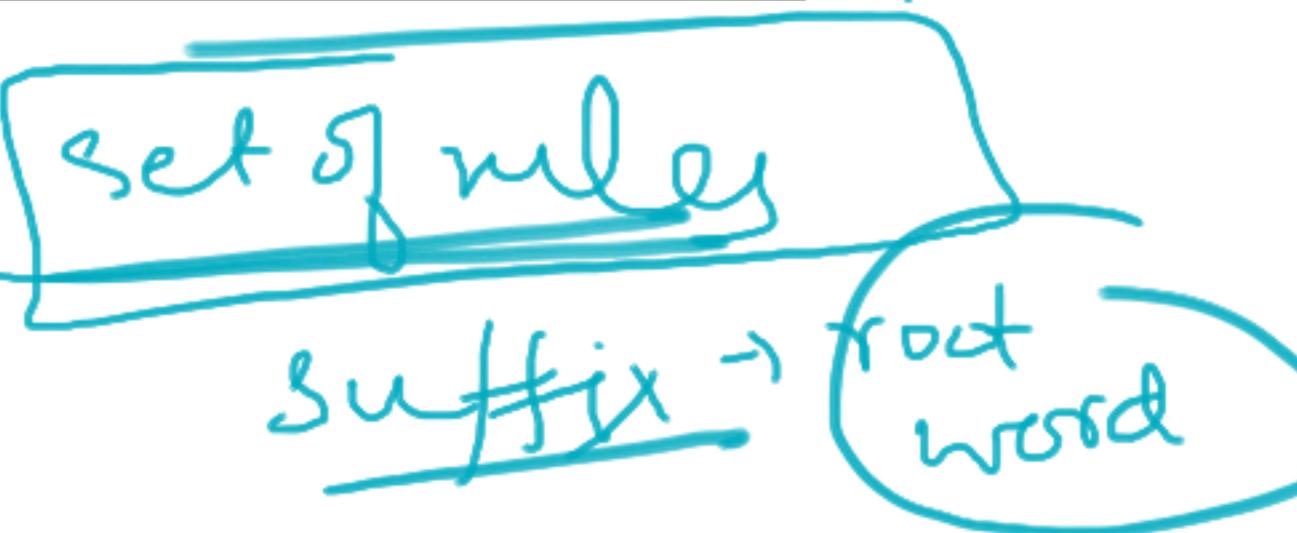
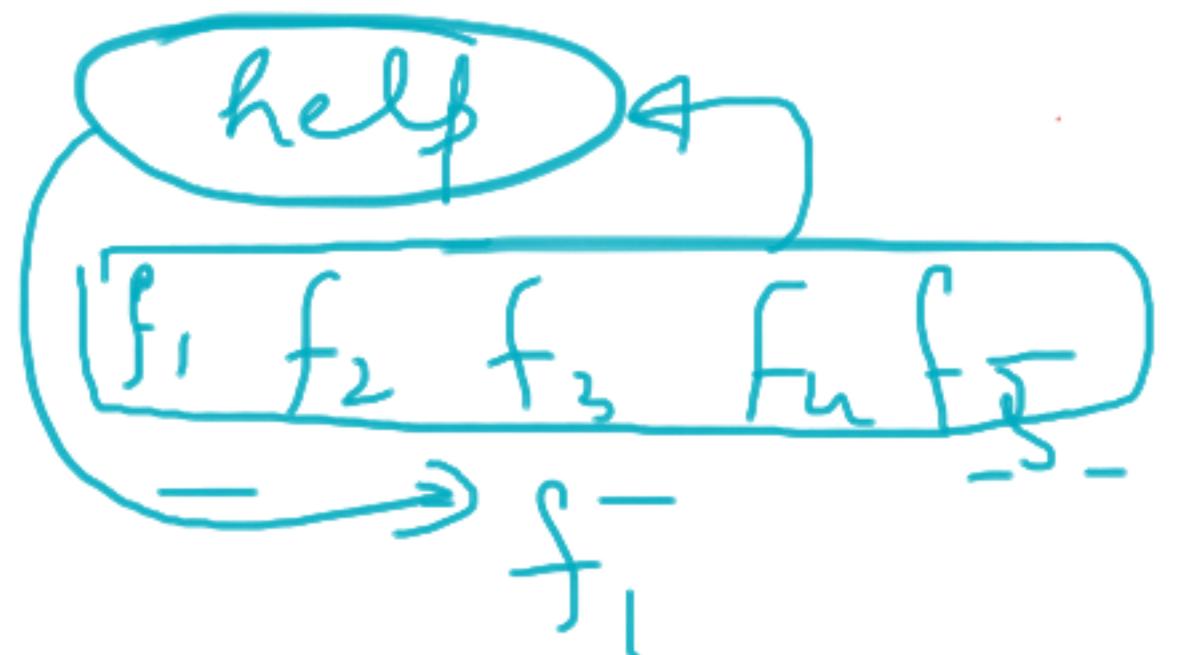
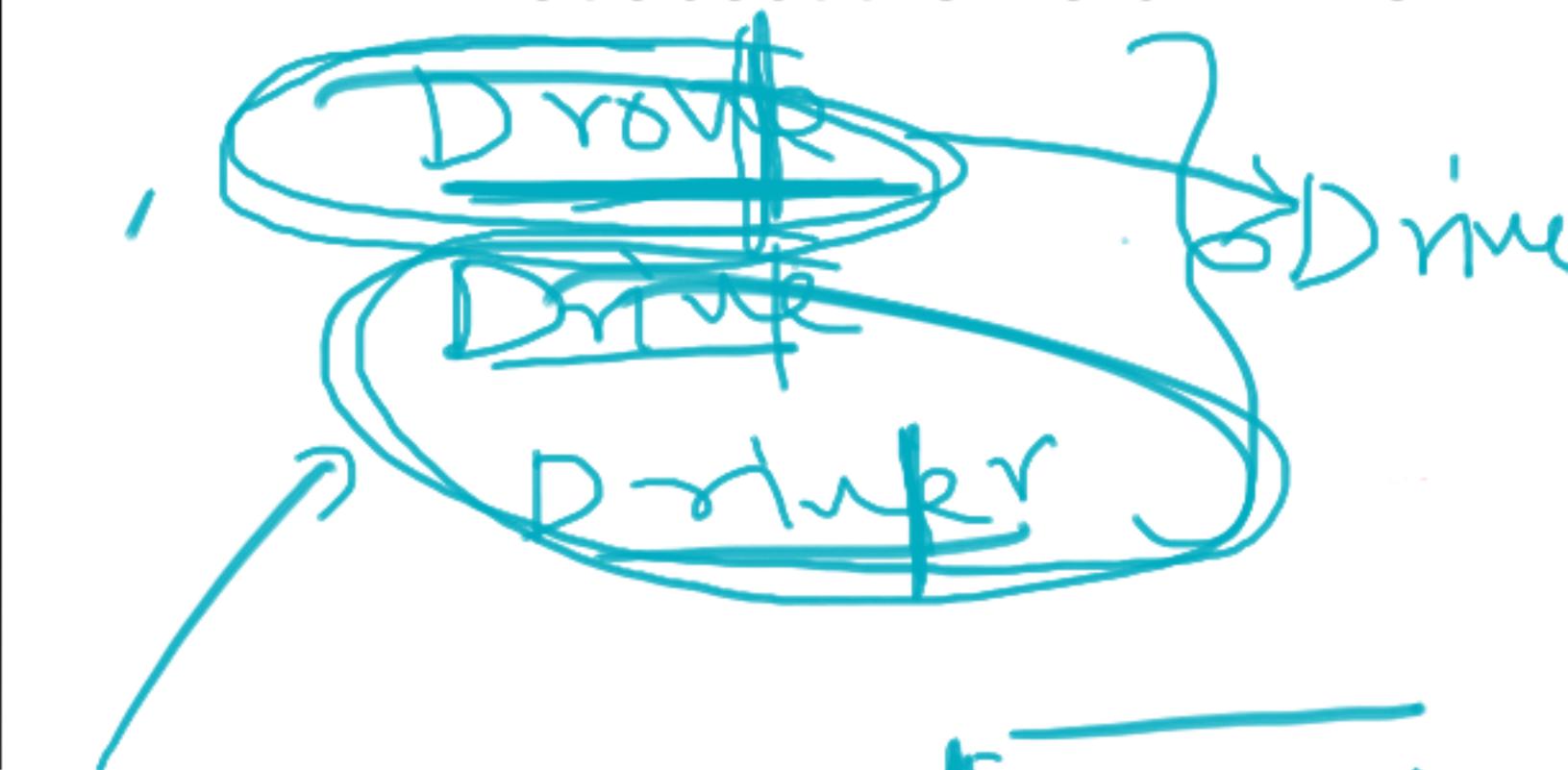
Form	Suffix	Stem
helps	-es	help
helping	-ing	help
helped	-ed	help
help	-	help
helper	-er	help

Removal of affixes

STEMMING



What about the word "Drive"



~~drive~~

Stemming algorithms

- Porter Stemmer: Most common and gentle stemmer.
- Fast, imprecise.

quick & dirty

Stemming algorithms

- Lancaster Stemmer (Paice-Husk): AGGRESSIVE. ITERATIVE.
- Simple, but heavy stemming can cause over-stemming.
- Leads to non-linguistic, with no meaning.



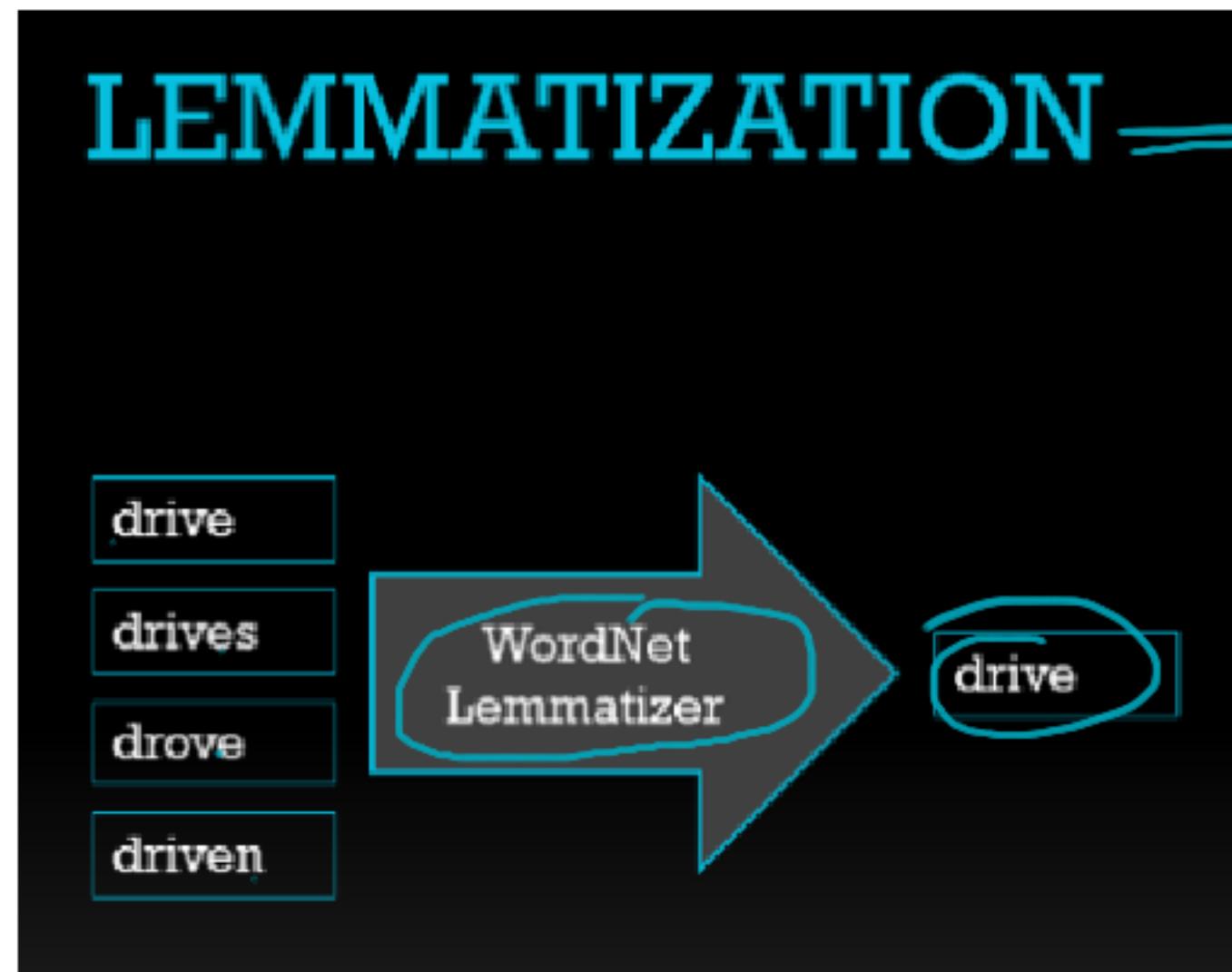
Stemming algorithms

- Snowball Stemmer (Porter2 / English stemmer): precision targeted over large datasets.
- Supports custom rules for any language.

imPsoe
widet

LS

Lemmatization



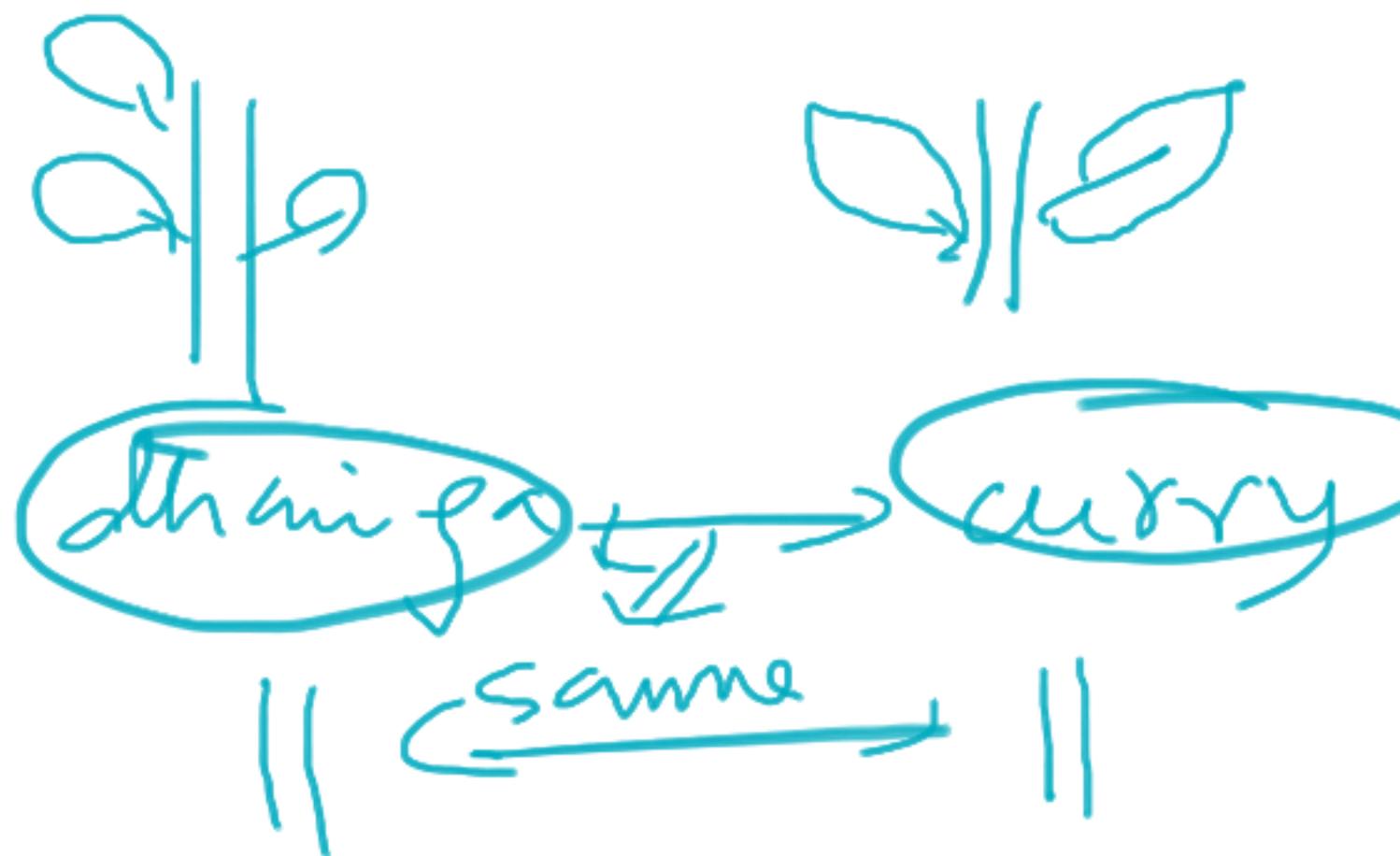
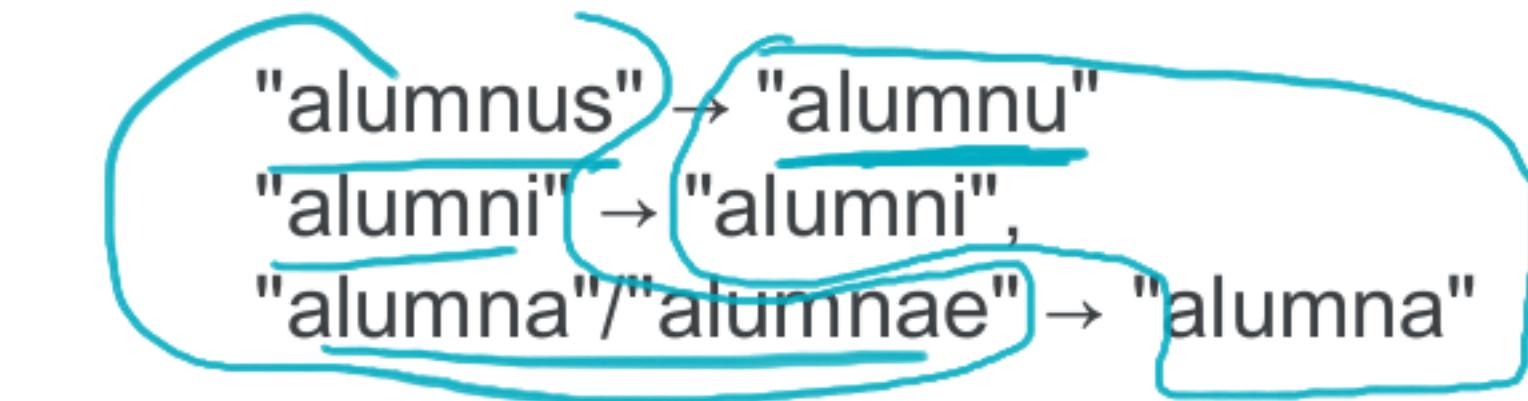
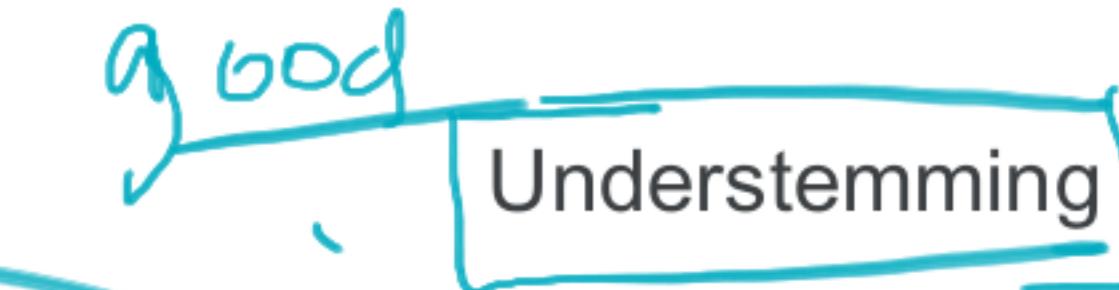
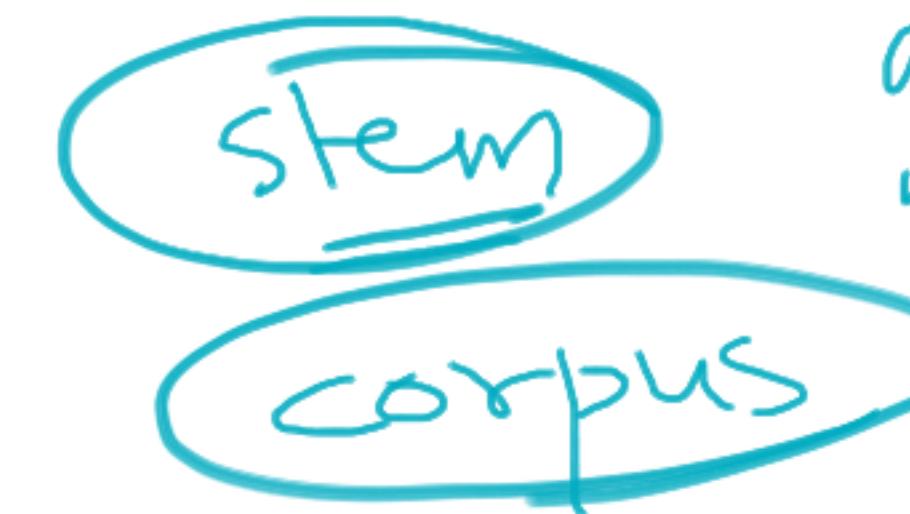
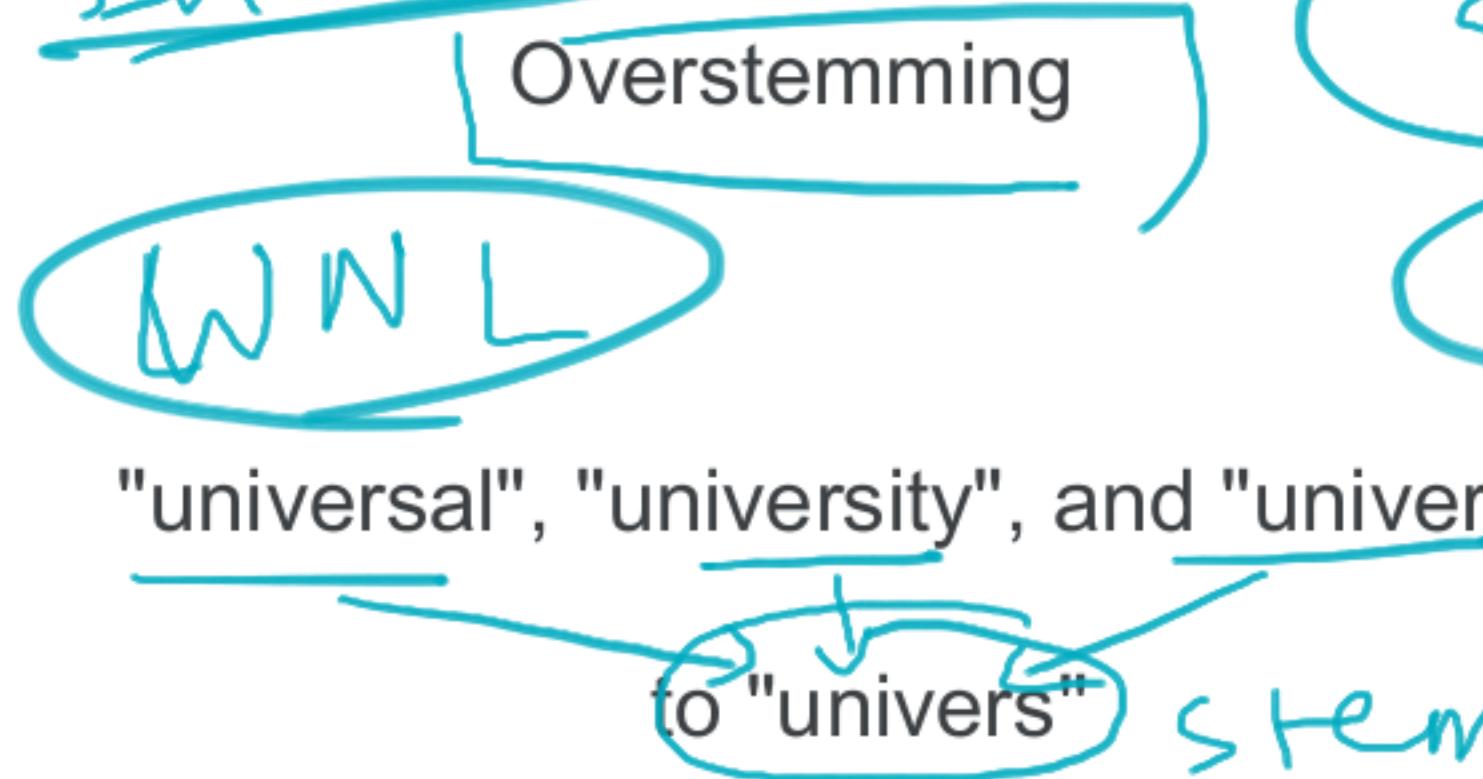
complex ~~root~~
dictionary

lot of time → processing high

Stem → Stemmed

Lemma → root meaning

Porterian



new novel

simple

Original	Stemmed	Lemmatized
visibilities	visibl	visibility
adhere	adher	adhere
adhesion	adhes	adhesion
appendicitis	append	appendicitis
oxen	oxen	ox
indices	indic	index
swum	swum	swim

ML model

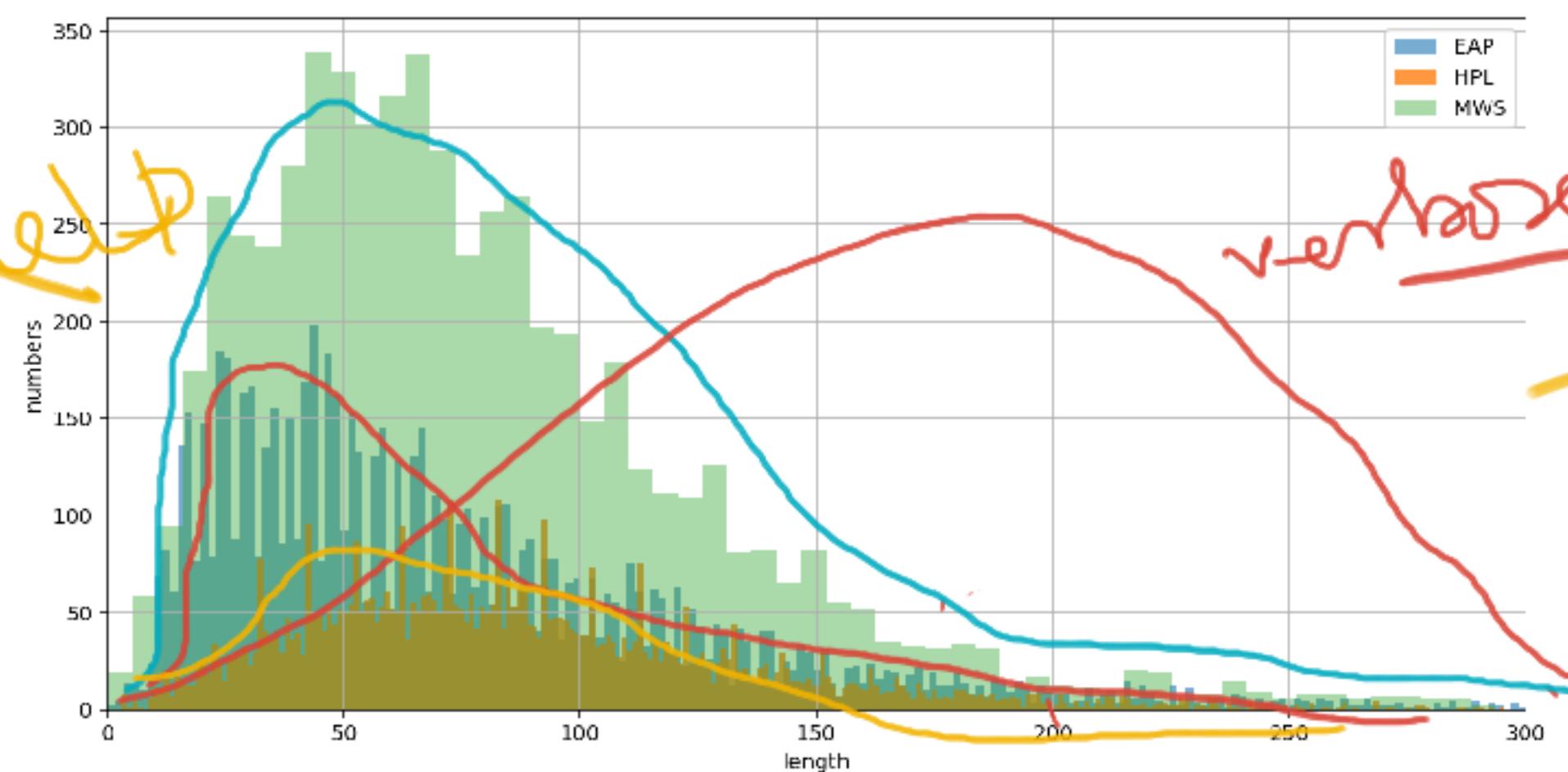
Advantages

few steps

Disadvantages

verbose

same dist



$$\underline{D+M} = ? \rightarrow$$

book

help

A hand-drawn diagram illustrating a chromatography process. At the top, three regions are labeled: 'M' (Mobile phase) on the left, 'T' (Stationary phase) in the middle, and 'C' (Chromatogram) on the right. An arrow points from 'M' towards 'T'. The central part shows a vertical column with horizontal dashed lines representing the stationary phase. Blue circles of different sizes are attached to the column, representing the separation of components. A yellow wavy line at the bottom represents the baseline.

Count words | terms
in DTM

relevant

Metro/City/Town studied ~~se~~ → els

A hand-drawn diagram of a trapezoid. The top horizontal side is labeled with the letter 'T' in the center. The bottom horizontal side is labeled with the letter 'M' in the center. The left vertical side is labeled with the letter 'M' at the top and the letter 'T' near the bottom. The right vertical side is labeled with the letter 'T' at the top and the letter 'M' near the bottom.

get-dummies

Nos

help half

- Tauf

is dog learn
NLP is ch

A hand-drawn diagram showing a sequence of components connected by vertical lines and horizontal arrows. The components are labeled: C, C, Ch, Go, Lok, Team, NLP, and Tauf. The diagram shows connections between Ch and Go, Lok and Team, and Team and NLP.

book learn

20

25

3D

1. 2. 3. →

h help

\Rightarrow count vectorize

book = how many doc contain book

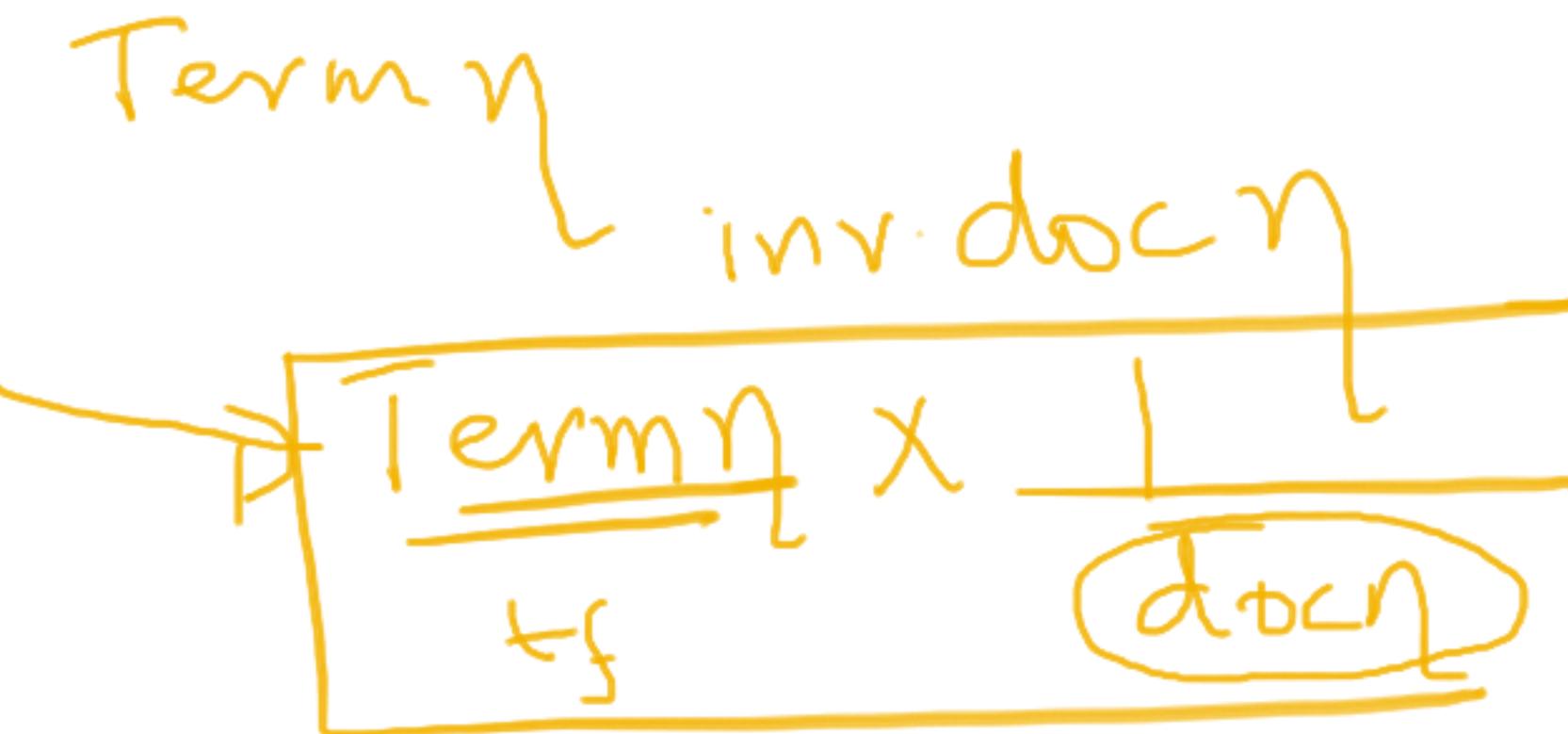
TF-IDF

d1: "sky blue"

d2: "sun bright today"

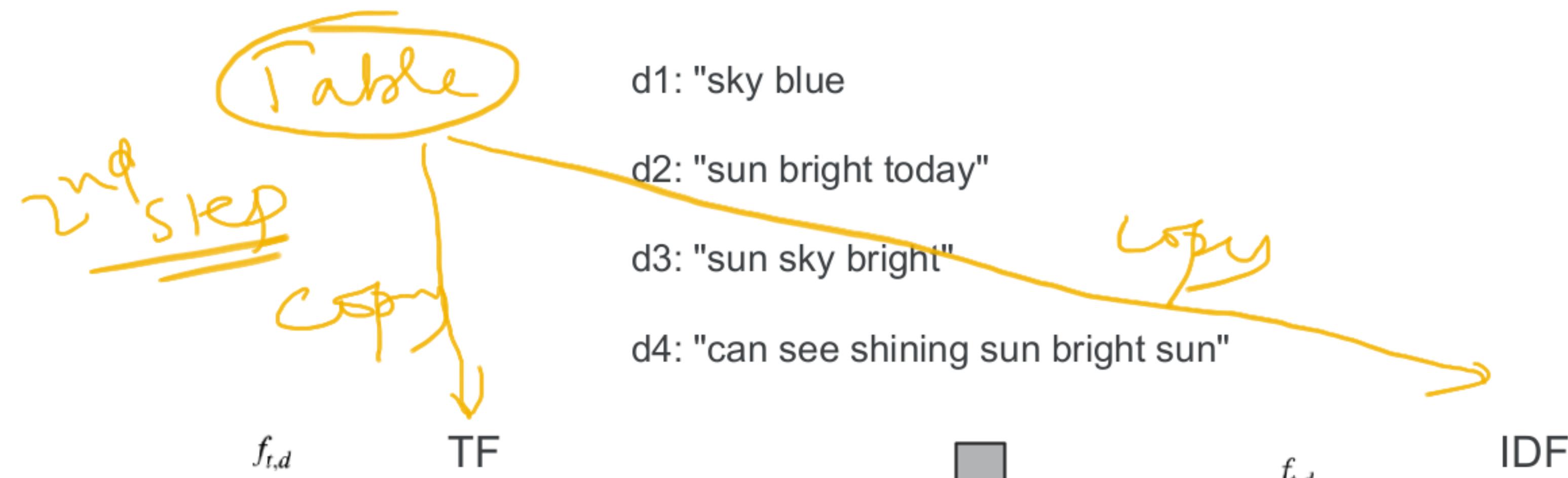
d3: "sun sky bright"

d4: "can see shining sun bright sun"



$f_{t,d}$ → ntable

	blue	bright	can	see	shining	sky	sun	today
→ 1	1	0	0	0	0	1	0	0
→ 2	0	0.1	0	0	0	0.1	0.1	0.1
→ 3	0	0.1	0	0	0.1	0.1	0	0
→ 4	0	0.1	0.1	0.1	0.1	0	0.2	0



	blue	bright	can	see	shining	sky	sun	today
1	1	0	0	0	1	0	0	0
2	0	1	0	0	0	1	1	1
3	0	1	0	0	0	1	1	0
4	0	1	1	1	1	0	2	0

	blue	bright	can	see	shining	sky	sun	today
1	1	0	0	0	0	1	0	0
2	0	1	0	0	0	0	1	1
3	0	1	0	0	0	0	1	1
4	0	1	1	1	1	1	0	2
n_t	1	3	1	1	1	1	2	3

Step-3 (i) - count one hot

$\text{idf} = \frac{n_t}{N}$ \rightarrow # of rows

$f_{t,d}$

$$\text{tf}(t, d) = \frac{f_{t,d}}{\sum_t f_{t,d}}$$

	blue	bright	can	see	shining	sky	sun	today
1	1	0	0	0	0	1	0	0
2	0	1	0	0	0	0	1	1
3	0	1	0	0	0	1	1	0
4	0	1	1	1	1	0	2	0

	blue	bright	can	see	shining	sky	sun	today
1	1/2	0	0	0	0	1/2	0	0
2	0	1/3	0	0	0	0	1/3	1/3
3	0	1/3	0	0	0	1/3	1/3	0
4	0	1/6	1/6	1/6	1/6	0	1/3	0

frac \Rightarrow step 3

$\text{tf}(t, d)$

	blue	bright	can	see	shining	sky	sun	today
1	1/2	0	0	0	0	1/2	0	0
2	0	1/3	0	0	0	0	1/3	1/3
3	0	1/3	0	0	0	1/3	1/3	0
4	0	1/6	1/6	1/6	1/6	0	1/3	0

	blue	bright	can	see	shining	sky	sun	today
1	1	0	0	0	0	1	0	0
2	0	1	0	0	0	0	1	1
3	0	1	0	0	0	1	1	0
4	0	1	1	1	1	1	0	0

$$\text{idf}(t, D) = \log_{10} \frac{N}{n_t}$$

$N = 10^2$ power log N

$\text{idf}(t, D)$

	blue	bright	can	see	shining	sky	sun	today
	0.602	0.125	0.602	0.602	0.602	0.301	0.125	0.602

- TF-IDF: Multiply TF and IDF scores, use to rank importance of words within documents

- Most important word for each document is highlighted

X

	blue	bright	can	see	shining	sky	sun	today
1	0.301	0	0	0	0	0.151	0	0
2	0	0.0417	0	0	0	0	0.0417	0.201
3	0	0.0417	0	0	0	0.100	0.0417	0
4	0	0.0209	0.100	0.100	0.100	0	0.0417	0

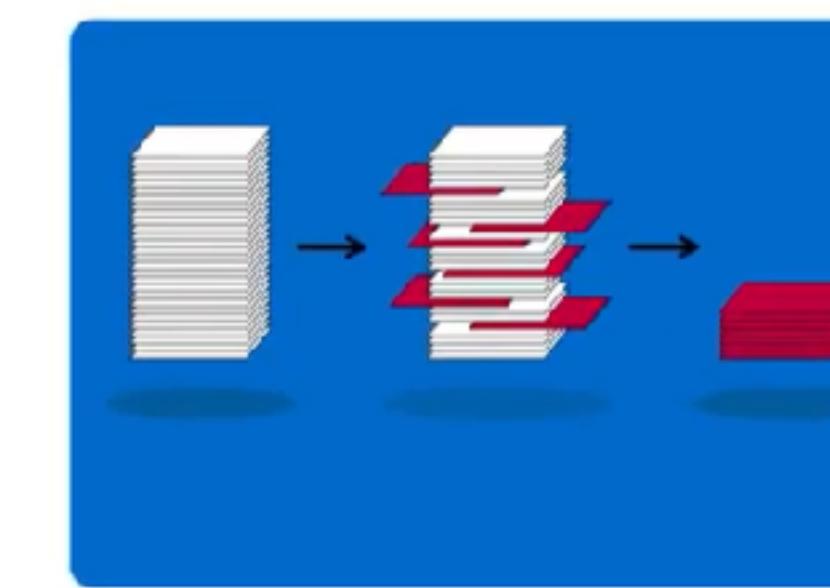
let's do it

imp

TF-IDF score computation. [Image Source]



It uses software that can identify concepts, patterns, topics, keywords, and so on in the data.

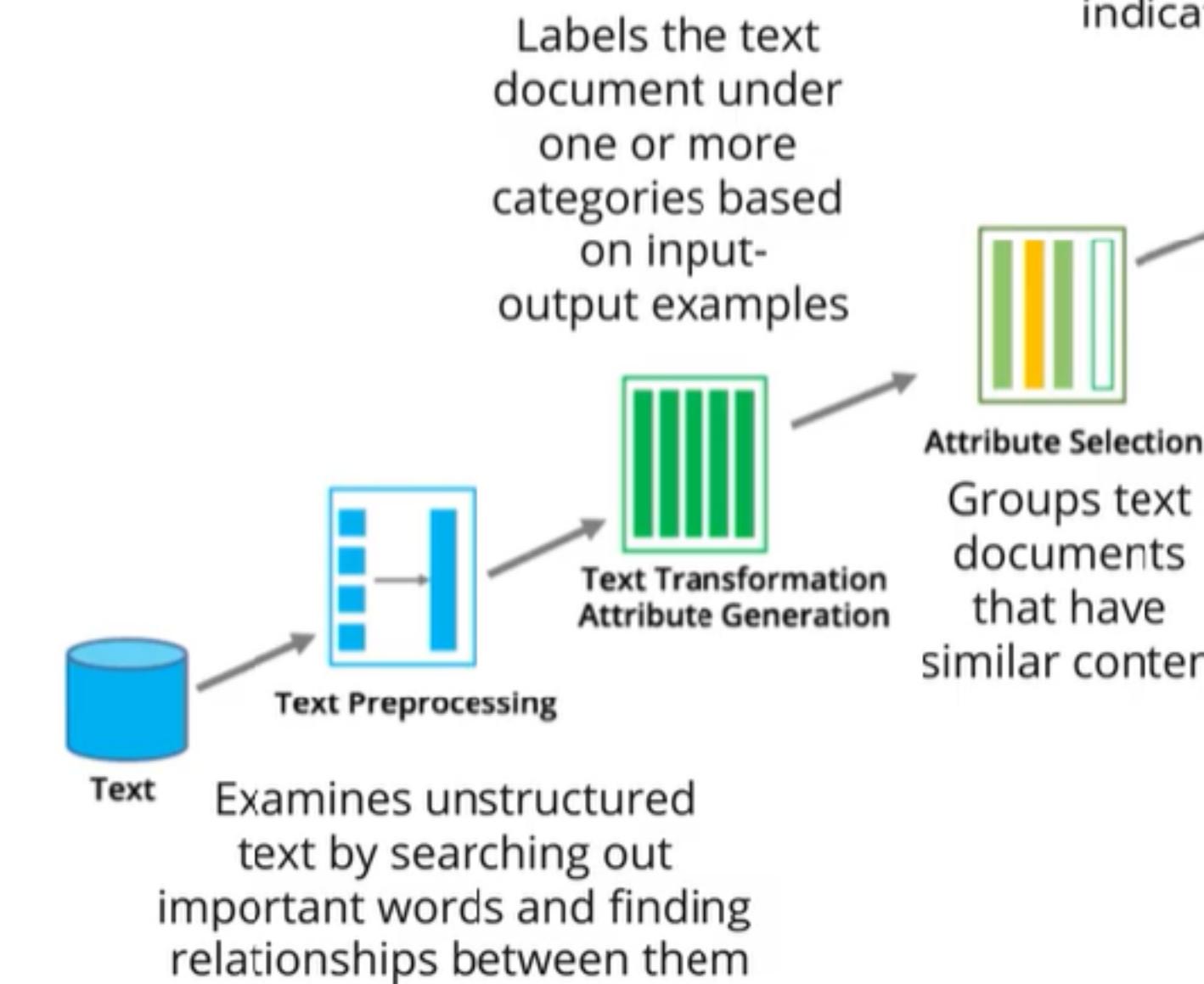


It uses computational techniques to extract high-quality information from unstructured text.

~~topic model~~

Flow of Text Mining

	id	text	author
0	id26305	This process, however, afforded me no means of...	EAP
1	id17569	It never once occurred to me that the fumbling...	HPL
2	id11008	In his left hand was a gold snuff box, from wh...	EAP
3	id27763	How lovely is spring As we looked from Windsor...	MWS
4	id12958	Finding nothing else, not even gold, the Super...	HPL
5	id22965	A youth passed in solitude, my best years spen...	MWS
6	id09674	The astronomer, perhaps, at this point, took r...	EAP
7	id13515	The surcingle hung in ribands from my body.	EAP
8	id19322	I knew that you could not say to yourself 'ste...	EAP
9	id00912	I confess that neither the structure of langua...	MWS



Uses text flags to represent documents and uses colors to indicate compactness



Interpretation or Evaluation

Reduces the length of the document by summarizing the details

Visualization



Attribute Selection

Text Transformation Attribute Generation



Groups text documents that have similar content

NLP Process Workflow

● TOKENIZATION

Splits text into pieces (tokens), remove punctuation.

STOPWORD REMOVAL

Removes commonly used words (Such as 'the') which are not relevant to analysis.

STEMMING AND LEMMATIZATION

● Reduces words to base form to be analyzed as a single item.

P.O.S TAGGING

Tags words to be part of speech (Such as verb, noun) based on definition and context.

● INFORMATION RETRIEVAL

Extracts relevant information from source