



POWERED BY

AI Outsider

Section 3

Text Representation

Want more? Head over to www.AIOutsider.com





Section Overview



Text Mining & NLP
Sentiment Analysis
Google Colab
Dataset Overview
Descriptive Statistics

Text Normalization
Features Cleaning
Tokenization
Stemming
Lemmatization

Text Representation
Negative/Positive
Bag-of-Words
TF-IDF





AIOutsider.com

WHY REPRESENTING?

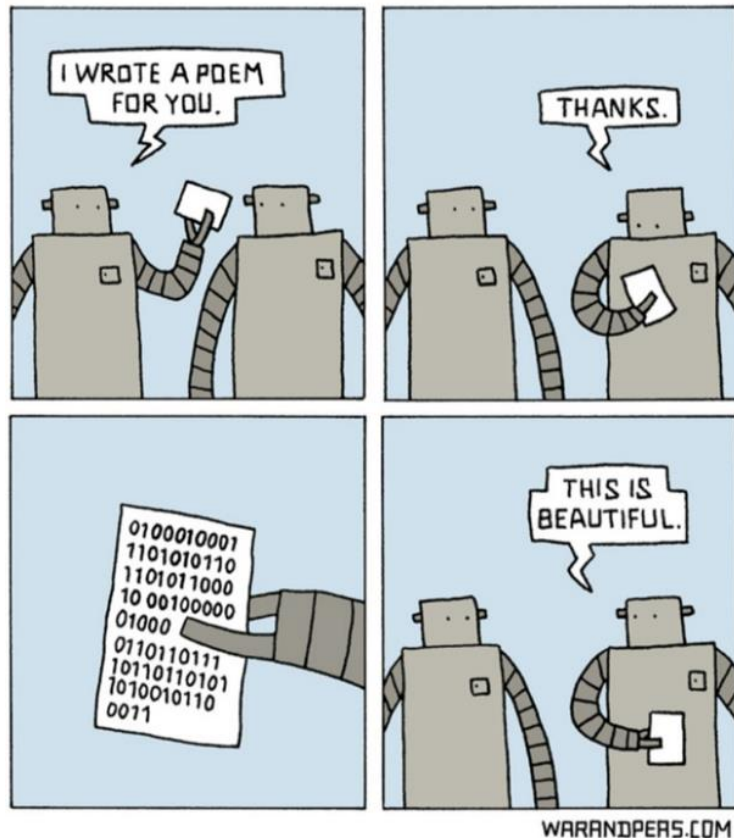


WHY REPRESENTING TEXT?



Representing text numerically allows it to be understood by ML models

Definition



AaI



10101
01011
10101

- ▶ Positive/Negative Frequencies
- ▶ Bag of Words
- ▶ TF-IDF





AIOutsider.com

POSITIVE/NEGATIVE



POSITIVE AND NEGATIVE FREQUENCIES



TWEET + « I am glad I got hired »

TWEET + « This is great »

TWEET - « This is bad »

TWEET - « I am sad I got fired »

	Freq(w, 1)	Freq(w, 0)
I	2 ←	2 ←
am	1	1
glad	1	0
got	1 ←	1 ←
hired	1	0
fired	0 ←	1 ←
this	1 ←	1 ←
is	1 ←	1 ←
great	1	0
bad	0 ←	1 ←



NEW TWEET « I got fired, this is bad »

$$X = \begin{bmatrix} \text{Pos Freq} , \text{Neg Freq} \end{bmatrix}$$

$$X = \begin{bmatrix} \Sigma \text{Freq}(w, 1) , \Sigma \text{Freq}(w, 0) \end{bmatrix}$$

$$X = \begin{bmatrix} 5 , 7 \end{bmatrix}$$





AIOutsider.com

BAG OF WORDS



➤ _ BAG OF WORDS REPRESENTATION

TWEET 1

« I like to learn »

TWEET 2

« We all like and all want to learn »



« to » « want »
« I » « and »
« like » « learn »
« all » « we »

words	I	all	we	like	and	want	learn	to
features	1	2	3	4	5	6	7	8

VECTOR 1

1	0	0	1	0	0	1	1
---	---	---	---	---	---	---	---

VECTOR 2

0	2	1	1	1	1	1	1
---	---	---	---	---	---	---	---

MATRIX DIMENSION = (# of tweets, # unique words in corpus)



AIOutsider.com

TF-IDF



TF

« TERM FREQUENCY »

IDF

« INVERSE DOCUMENT FREQUENCY »

$$tf_{w,d} = \frac{n_{w,d}}{\sum_k n_{w,d}}$$

$$idf_w = \log\left(\frac{N}{df_w}\right)$$

TF



IDF



TF-IDF

features

TWEET 1

« I like my cat »

TWEET 2

« I love my dog »



	d1	d2		d1	d2
w ₁	1/4	1/4	w ₁	Log(2/2) = 0	0
w ₂	1/4	0	w ₂	Log(2/1) = 0.3	0.075
w ₃	0	1/4	w ₃	Log(2/1) = 0.3	0
w ₄	1/4	1/4	w ₄	Log(2/2) = 0	0.075
w ₅	1/4	0	w ₅	Log(2/1) = 0.3	0
w ₆	0	1/4	w ₆	Log(2/1) = 0.3	0

: \> _ FOLLOW US!



AIOutsider.com



twitter.com/AIOutsider_

Disclaimer

All Trademarks referred to are the property of their respective owners.

© 2021 – All rights reserved.

