

→ Text Classification:

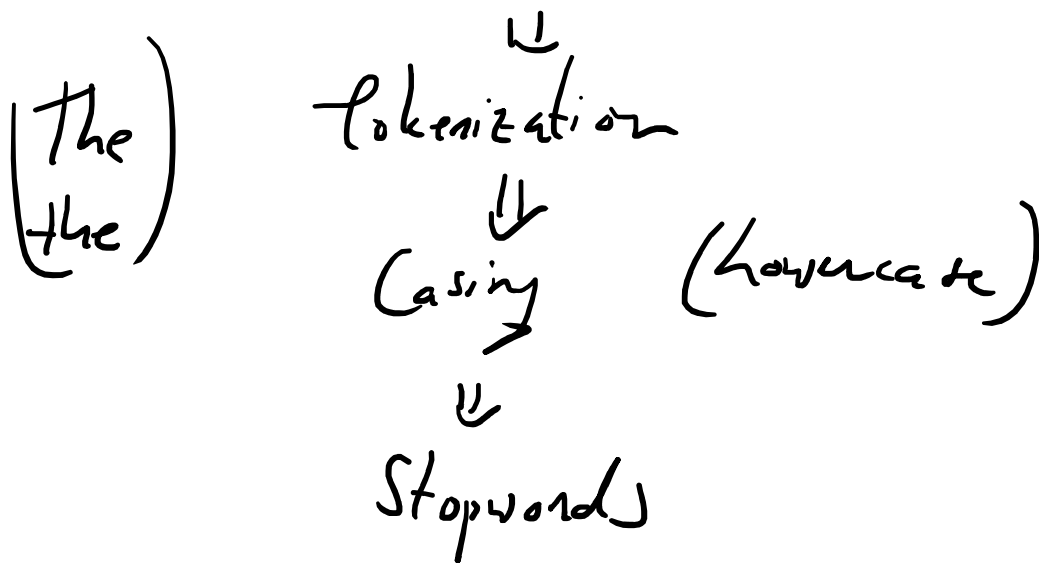
NLP	ML
-----	----

→ Discrete Samples

[0, 1, 2, 3]

→ Preprocessing
 ↓
 Embedding
 ↓
 Model Building

→ Character Process (Regex)
 ↓
 Effective Splitting (Sentence Conversion)



⇒ Embeddings:

↳ BoV: Bag of Words

"I like to go to the movies"

"I do not like movies like this."

↳ Uses Stochastic Sparse Matrix.

	BoV ₁	BoV ₂	TF-IDF ₁	TF-IDF ₂
I	1	1		
like	1	2		2/7 lg 7
to	2	0		
go	1	0		
the	1	0		

movies	1	1
do	0	1
not	0	1
this	0	1

→ TF-IDF: Term Frequency - Inverse Document Frequency.

$$TF(t, d) = \frac{\text{Occurrence of } t \text{ with } d}{\text{Total words in } d}$$

Document
All docs
count(D)

Term
Term

$$IDF(t, D) = \log \frac{N}{\text{documents with } t}$$

$$TF-IDF(t, d, D) = TF(t, d) \times IDF(t, D)$$