

Raw Data → Cleaning → Normalization

Python libraries &

NLPK (slow)

Sparcy (fast)

1) Education, Research, Experimentation

1) Industry & Production.

2) Modular & provides a lot of tools & resources for nlp tasks.

2) Pipeline based & provides fast & integrated api for nlp tasks.

- Modular but basic building blocks for nlp tasks.

- Provides limited but highly accurate nlp tools & resources.

3) Good for analytics & exploration.

3) Stable & production friendly.

4) Large community.

4) Growing community.

5) Requires more manual tasks.

5) Less manual work.

Text Cleaning →

1) Remove unwanted characters →

2) `str.translate(chars to replace, chars to replace, with chars to delete)`

- Creates a mapping for characters & their replacement

3) `text.translate(mapping from make-trans)`

- Does the actual replacement in text.
- Can use a custom dictionary instead of the one provided by `make-trans`.

ex 2 `text.translate(str.translate('!', ',', str.punctuation))`
- removes punctuations from text.

2) Lowercase → `sent.lower()`

3) Remove url, whitespace & digits →

`import re`

Pattern for url.

`pattern = re.compile('or 'http s?://\S+ | www. \S+')`

`pattern.sub(repl to replace with, original text)`

- Replace all the matched text with replacement text.

4) whitespace →

`re.sub(r'\s+', ' ')` → Replace extra whitespaces with a single white space.

`re.sub(r'\t+', '')`

Tokenization → Breaking text into smaller sections, called tokens.

`with tokenizer`
`from spacy import管nt_tokenizer, word_tokenizer.`
`import emoji`

`sentences = sent_tokenizer(raw_text)`

`for i, sentence in enumerate(sentences):`

`print(f'Sentence {i+1}: {sentence}')`

`nlp = spacy.load('en_ner_web-sm')` → Loads the spacy pipeline.
`spacy.set = nlp(raw_text)`

- Spacy passes the entire raw text through its pipeline.

Stop words → Words that occur frequently.

- This, in, a, the...

- Don't have much value but reduces the worthiness of valuable words.

from with . removes import stopwords

words = word_tokenizer(text)

words2 = [word for word in words if word not in set(stopwords.words('english'))]

— using spacy &

words_spacy = [word.lemmatize() for word in spacy.Doc if not word.is_stop]

→ Stemming & Lemmatization → (both deal with bringing the word to its root form)

↳ applies to bring words to root form.

↳ by stripping characters at the end of the words.

↳ fast but can be inaccurate.

↳ applies vocabulary & morphological analysis to bring words to their root form.

↳ slow but accurate.

with & PorterStemmer(), WordNetLemmatizer.

stemmer = PorterStemmer()

stem_words = [stemmer.stem(word) for word in words2]

// using spacy &

stem_words = [word.lemma_ for word in spacy.Doc if not word.is_stop]

Sentiment Analysis →

genism → Python library for gen to -
 polarity & Sentiment [-ve < 0 < +ve]
 Subjectivity & opinion are a part.

1) Tint Blob → Library for sentiment & subjectivity analysis.

from tintblob import TintBlob

blob = TintBlob(text)

blob.sentiment.polarity, blob.sentiment.subjectivity.

& polarity → Information about sentiment of the sent.

→ Subjectivity → If the sent is subjective or objective
 ↴ opinonated part

(objective) ① ————— | (subjective)

2) NRC lexicon → Maps words to 8 basic words that
 convey the different emotions of the sent.

ME words
 manually made

- ultimately maps to two emotions.

- Can analyse multiple emotions beyond good & bad.

sent → 8 emotions.

→ 2 sentiment.

from nrclex import NRCLex

emotion = NRCLex(text)

emotion.words

- raw emotion scores → freq. of occurrence of each word!
- top emotions
- affect frequencies → freq. of each score/total score want.

3) Affin → Simplified sentiment analysis procedure.

2500 words.

- fast & light-weight.

- provides a score b/w -5 to 5.
- aggregates the score for each word for a sentence.

from affin import Affin

af = Affin()

score = af.score(text)

af.split(text) → Split text into words.
af.findall(text) → return all the emotional text.

.score_with_pattern(text) → return score of each emotion.
.score_with_wordlist(text) → total score of emotions.