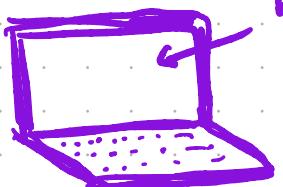
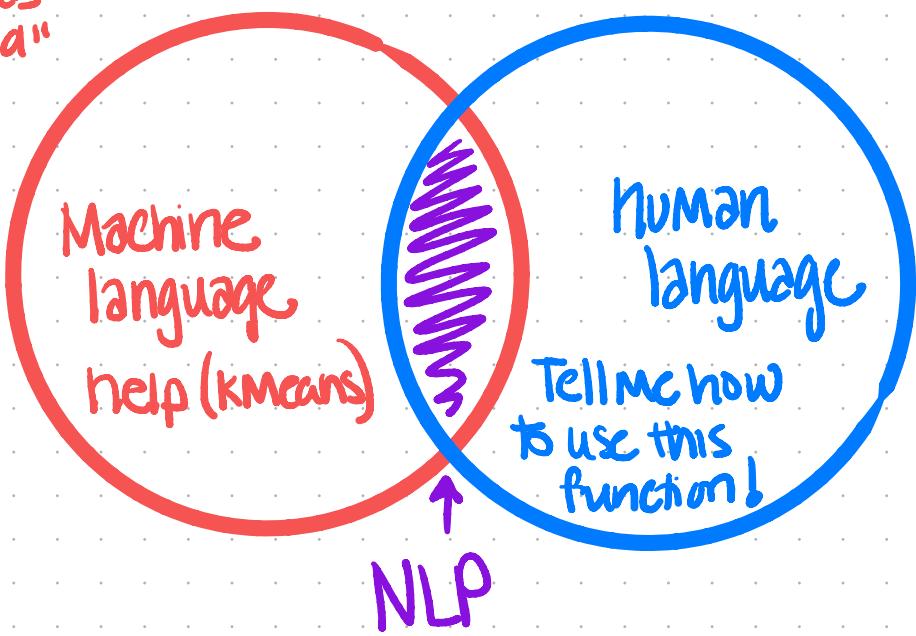


NLP

Natural
Language
Processing

Processing & analyzing large amounts
of natural language data via
programming.



USE CASES for NLP

Voice of the Customer

Analysis of surveys, tweets, other feedback
What drives + / - responses? What is being asked for?
What drives engagement?

Semantic Search

a search engine that searches using intent & context rather than just key words.

Knowledge Management & Discovery

tag documents, blog posts, etc. based on the topics for ease of recommending, finding, discovering.

Customer Support

Route support requests based on needs / topic identified through the text contained in the request.

Chat bots to solve the common issues before routing to a human.

Security

Phishing emails, national security / terrorism on Social Media, dark web
cyberbully

Healthcare

New cells of viruses/flu through SM monitoring, Mental health via SM

Virtual Assistants

Siri, Alexa, google home

Translation, Speech to text, competitive intel, Spam

Spell check & Autocorrect, Auto completion,

NLP Module Goals & Agenda

- USE **REGULAR EXPRESSIONS** to find/extract/replace/substitute text that matches a pattern.
- Acquire data through **WEB SCRAPING**
- Parse text : **NORMALIZE, TOKENIZE, STEM, LEMMATIZE, remove STOPWORDS**
- Explore text data **WORD CLOUDS, N-GRAMS**
- Create features using **BAG OF WORDS, TF-IDF**
- **PROJECT**

TYPES of NLP PROGRAMS

Sentiment Analysis



Topic Modeling

"Data Science is fun"

→ No label - extract topics (LDA)

→ tag = DataScience LABEL!

Text Classification

Translation

Google's Crowdsource App → Get labeled data
NSO, Some Rule-based using dictionaries

Word Embeddings

find similar articles, forum question duplication

Document Embeddings

Supervised	unsupervised
✓	✓
✗	✓
✓	✗
✓	✓
✓	✓
✗	✓

WRANGLING in NLP

MODULES

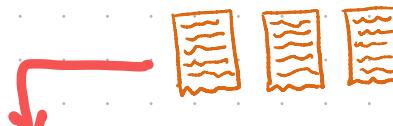
Acquire your **CORPUS** of documents



REQUESTS
(get the HTML)

BEAUTIFUL
SOUP
(parse the HTML)

Convert your corpus to individual **DOCUMENTS**

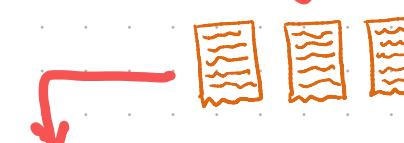


UNICODE DATA
RE

NORMALIZE text

He went to the store to buy ~~X~~ rolls of toilet paper ~~X~~. He was really
~~nave~~ to think there would be toilet paper.

he went to the store to buy rolls of toilet paper he was really
nave to think there would be toilet paper



NLTK.tokenize.toktok

Break documents into **TOKENS**



NLTK.porter
NLTK.stem

Create **STEMS** or **LEMMAS**

went → go rolls → roll was → is

NLTK.stopwords

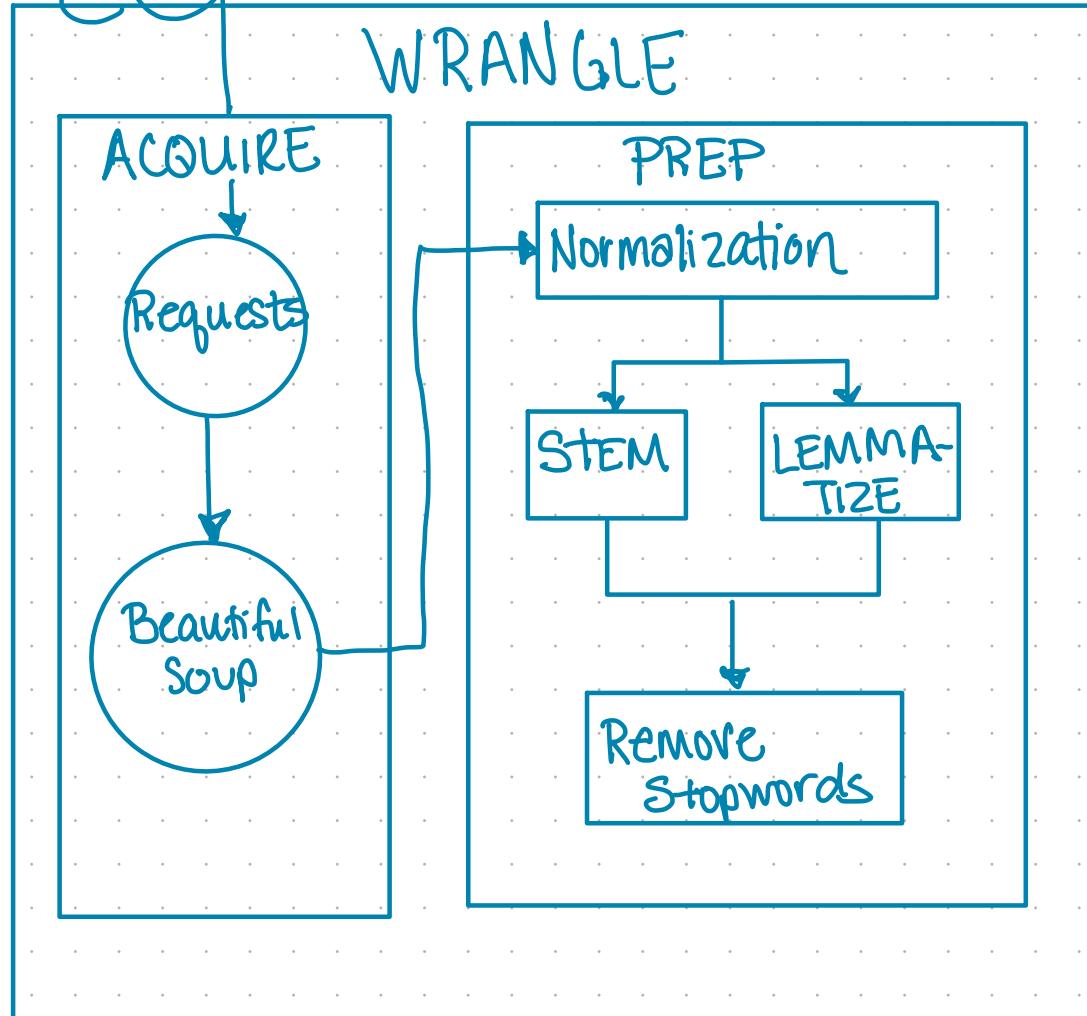
Remove **STOPWORDS**



EXPLORE

NLP DATA FLOW

ex.com



EXPLORE

term-frequencies*

word clouds

n-grams*

MODEL

Tf-idf

term-freq. -
inv. doc. freq.

Classifier

Bag of
Words

NLP METHODS

PROS

Rule-based

Explicit Logic, Human Readable

Statistical
(e.g. Bayesian)

Automatable, Easy to develop

Vector Methods
(e.g. Word2vec)

Uses common and easily
accessible ML algorithms

Deep Learning
/ Neural Nets

Good results, little manual
effort.

CONS

Time to dev, maintain, doesn't
scale well.

Doesn't generalize well & not
flexible, doesn't capture semantics

Doesn't capture document
structure, or the "Big Picture"
well.

Black Box, Risk of discrimination,
biased decisions made as
result, unexplainable,
high computational cost.

NLP Vocabulary

ENTITY any word or series of words that consistently refers to the same thing.

CORPUS body of docs, entire sample using to analyze

DOCUMENT - single observation

STEM - removing prefixes, suffixes, to the core string. (easing → eas)

LEMMA - the base word. (easing → ease)

TOKEN - single linguistic unit

CORPUS → **DOCUMENT** → **TOKEN**