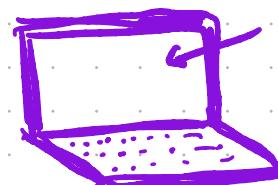
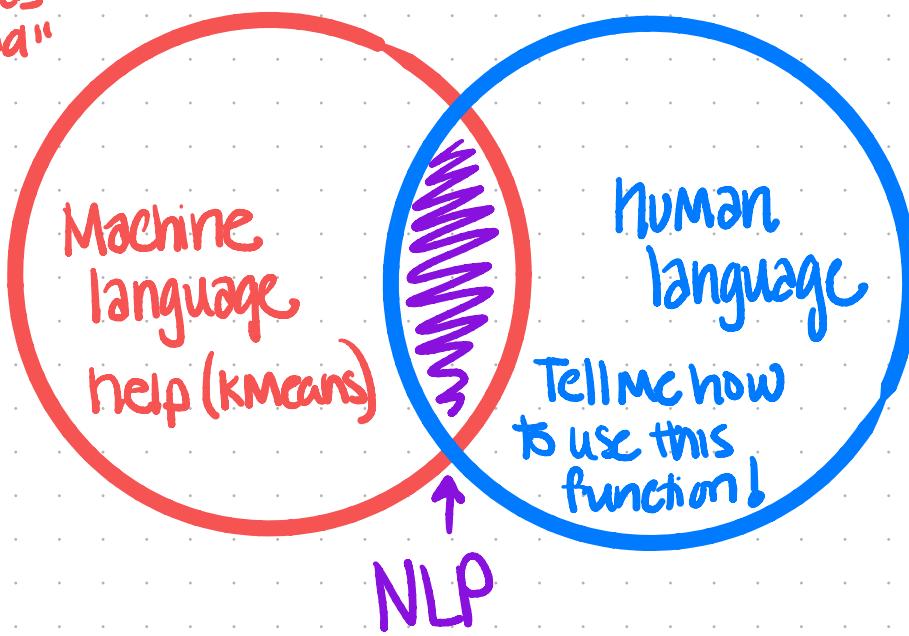


NLP

Natural
Language
Processing

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

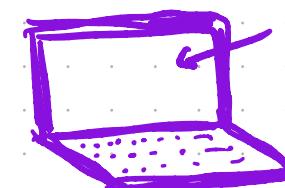
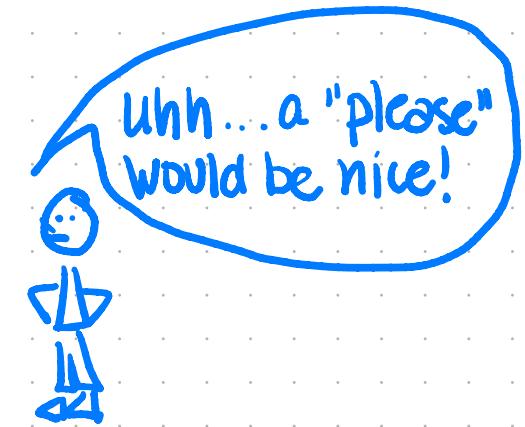
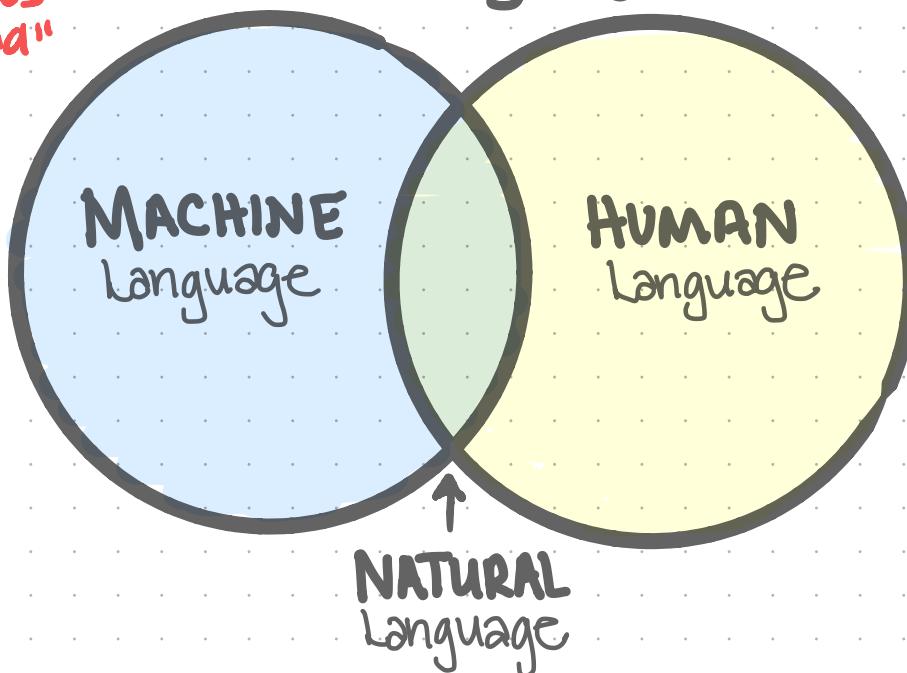


NLP

Natural
Language
Processing

{ Processing = analyzing large amounts of
natural language data via programming }

The Venn Diagram
of Language



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,

TYPES of NLP PROGRAMS

Sentiment Analysis



You're the s**t!
You're full of s**t!
My dog takes a s**t
on our walks.

Topic Modeling

"Data Science is fun"

→ No label - extract topics (LDA)

Text Classification

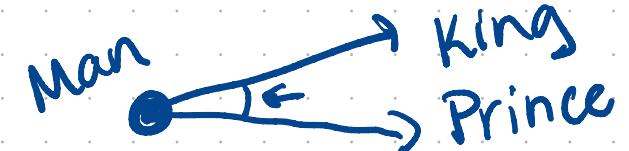
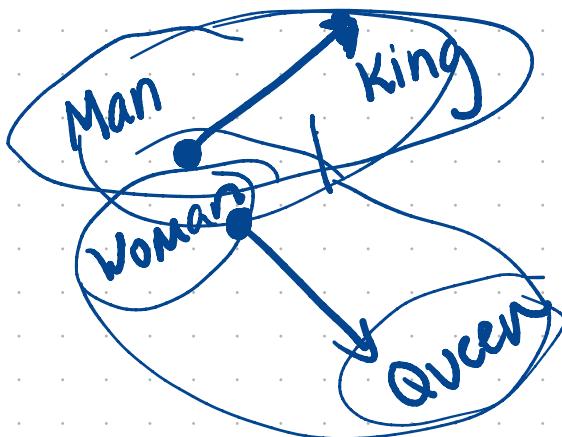
→ tag = DataScience LABEL!

Translation

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

Word Embeddings →
Document Embeddings

find similar articles,
forum question duplication



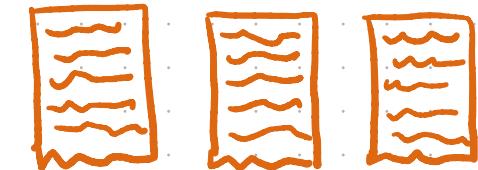
Supervised
Unsupervised

THE NLP PIPELINE - WRANGLE

Acquire your **CORPUS** of documents



Convert your corpus to individual **DOCUMENTS**



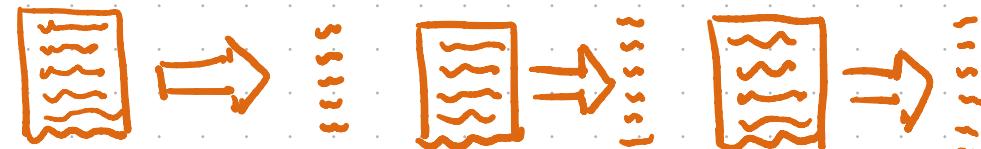
NORMALIZE text

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



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

Break documents into **TOKENS**



went → go rolls → roll was → is

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

Remove **STOPWORDS**



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 -

CORPUS - Body of docs, entire sample using to analyze

DOCUMENT - single observation

STEM - removing prefixes, suffixes - to the core string.

LEMMA - the base word.

TOKEN - single linguistic unit

CORPUS → **DOCUMENT** → **TOKEN**

NLP Stemming vs. Lemmatization

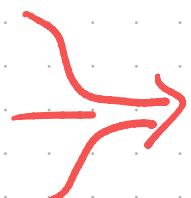
Stemming: ~~apply~~ helpful asked playing toys
↓ ↓ ↓ ↓ ↓
apply help ask play toy

Stemming is fast & simple to use. BUT:

news caring tried
↓ ↓ ↓
New car tri

Lemmatization: uses dictionaries to reduce a word to its base word
in order to group all words with the same base together.
It requires more computational resources than stemming.

Caring
Care
Cared



tried
trying
retry



NLP MODELING TECHNIQUES

LDA - Latent Dirichlet Allocation - Topic Modeling, Unsupervised

Data Science is a growing field. 100% Data Science - (Topic A)

Web developers build applications 100% Web Dev (Topic B)

Codeup offers a web development and data science program.

50% Data Science (Topic A)

30% Web Dev (Topic B)