# Text analysis

≈ Natural Language Processing, or "*How to do cool stuff with words.*"



**Emily Rae Sabo**   Data Camp | June 19, 2019

# 2 objectives for this session:

- ✓ What is NLP /Text Analysis and why would I use it?

- ✓ What tools are out there for me to use?
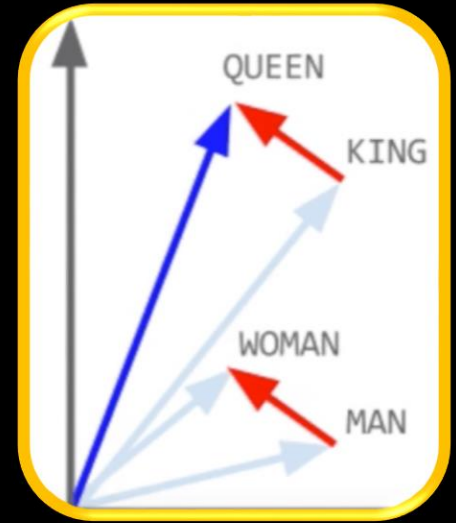
# What is NLP used for?



**Predicting language**

**+**

**Translating language**

**+**

**Finding patterns in language**

**+**

**Measuring meaning in language**

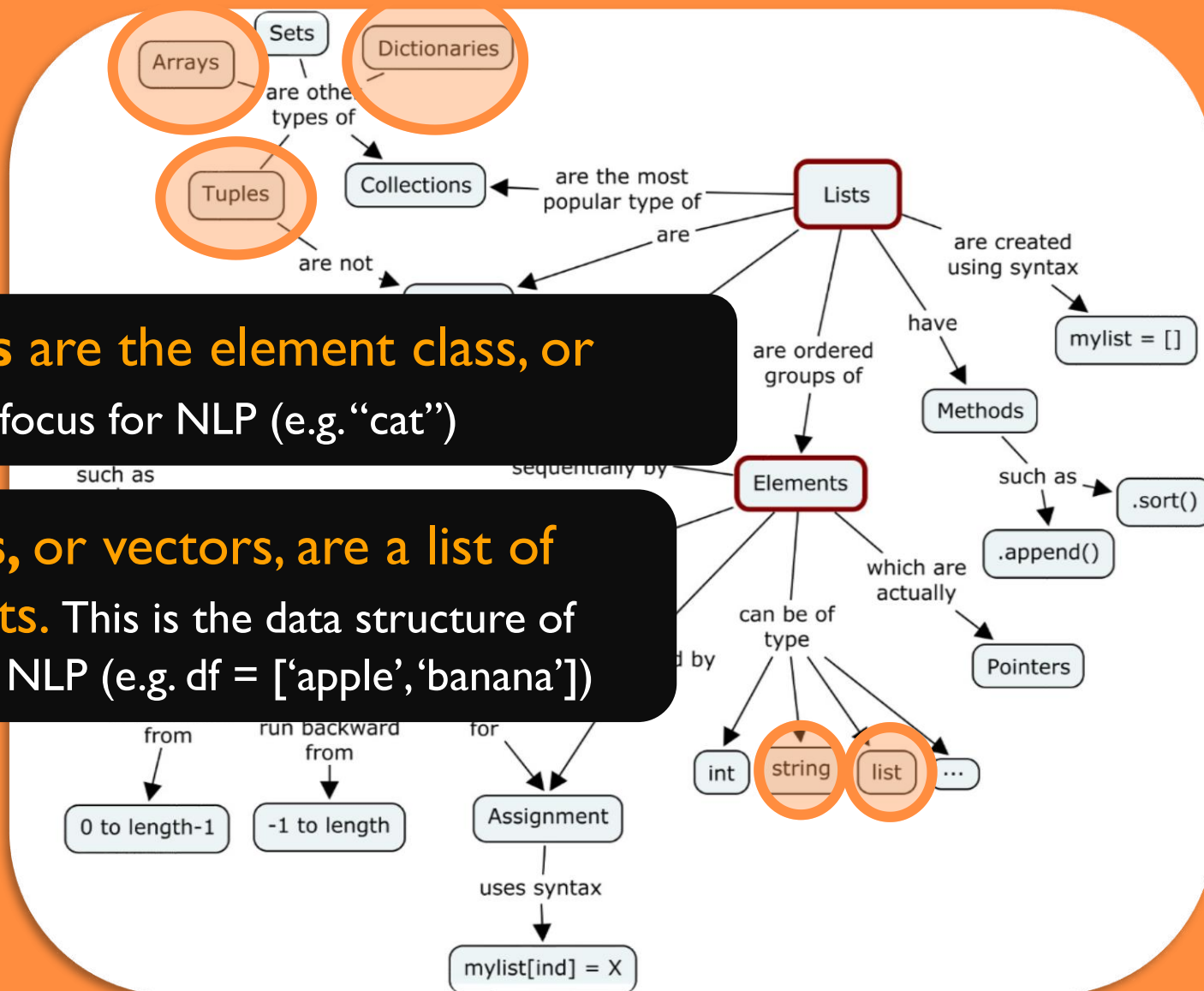# How to apply Text Analysis

Finding patterns in language



Measuring meaning in language



- Change over time with Google Ngram
- Topic Modeling with Gensim, NLTK
- String matching and token extraction with RegEx

- Vector space modeling with word-embedded vectors like Word2Vec in *Gensim* or GloVe in *SpaCy*

# Python's basic elements & data structures



**Strings** are the element class, or type, of focus for NLP (e.g. "cat")

**Arrays,** or vectors, are a list of elements. This is the data structure of focus for NLP (e.g. df = ['apple', 'banana'])

Within the concept map:

- Sets
- Arrays
- Dictionaries
- Tuples — are other types of
- Collections ← are the most popular type of — Lists
- are not
- Lists are created using syntax — mylist = []
- Lists have Methods
- are ordered groups of
- such as
- sequentially by — Elements
- Methods such as → .sort()
- .append()
- which are actually → Pointers
- can be of type
- from → 0 to length-1
- run backward from → -1 to length
- for → Assignment
- int, string, list, ...
- Assignment uses syntax → mylist[ind] = X

# 4 TAKE-AWAYS

1. **Google Ngram Viewer** is a quick 'n dirty tool for measuring word frequency change over time.

2. **Topic modeling** is a dimensionality reduction technique used to reveal "topics" in a document.

3. **Regular Expressions (RegEx) is the syntax you use** to do string matching, text cleaning, and token extraction.

4. **Word-embedded vectors** are decomposed matrices from a huge word matrix that tells you about word meaning.

# How to measure changes in word frequency over time?

## Google Ngram Viewer

- The founding tool of "culturomics"
- Advantages vs. limitations?
- Share one way you could imagine using this in your research.
- Go and play!
  - https://books.google.com/ngrams
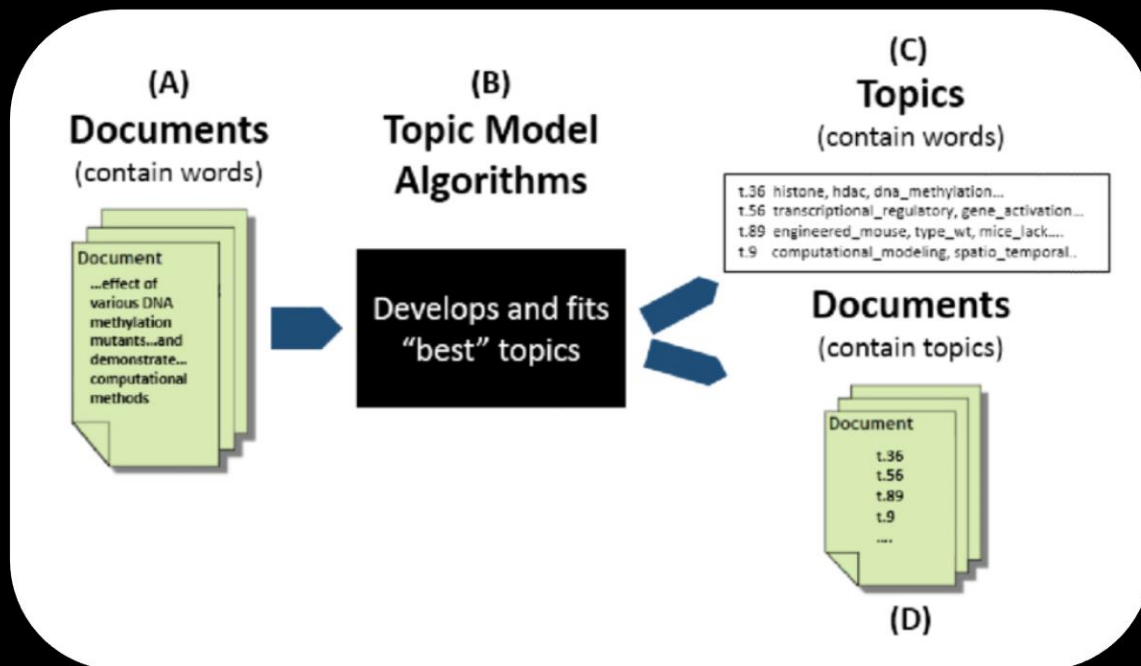  - https://books.google.com/ngrams/info

# What is Topic Modeling?

It is an **unsupervised approach** used for finding and observing the bunch of words (called "topics") in large clusters of texts."
*Bansal (2016)*

Click here for a good starter on Topic Modeling in Python with NLTK and Gensim

- It's a dimensionality reduction technique used to discover the hidden or abract "topics" that occur in a document or collection of documents.

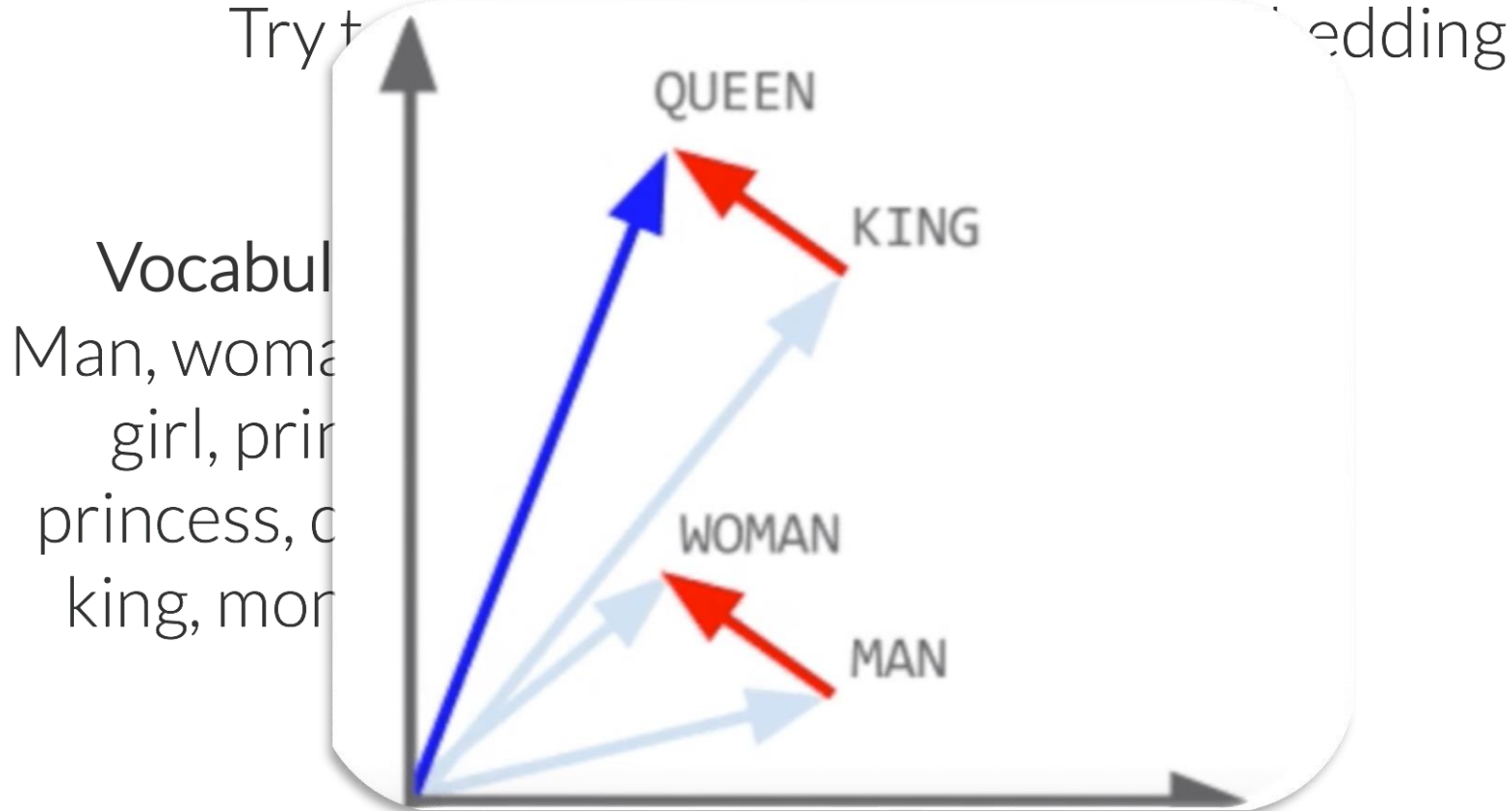- Techniques you may have heard of before: LSA (Latent Semantic Analysis) and LDA (Latent Dirichlet Allocation)



(A) **Documents** (contain words)

Document
...effect of various DNA methylation mutants...and demonstrate... computational methods

(B) **Topic Model Algorithms**

Develops and fits "best" topics

(C) **Topics** (contain words)

t.36 histone, hdac, dna_methylation...
t.56 transcriptional_regulatory, gene_activation...
t.89 engineered_mouse, type_wt, mice_lack....
t.9 computational_modeling, spatio_temporal..

**Documents** (contain topics)

Document
t.36
t.56
t.89
t.9
....

(D)

# What are Regular Expressions, or RegEx?

1. **Work through one tutorial:** https://regexone.com/ https://www.tutorialspoint.com/python/python_reg_expressions.htm

2. **Then, open Jupyter,** create your own mini-corpus (~20 words) and **write RegEx code to match a string from your corpus.**



**Pro-tip reminders:** Be computational *and* creative in your approach. There are an infinite number of ways to accomplish a string matching task!

# Vector Space Modeling, Word-embedded vectors & Cosine Similarity

Try t...                                    ...edding

Vocabul...

Man, woma...

girl, prir...

princess, c...

king, mor...
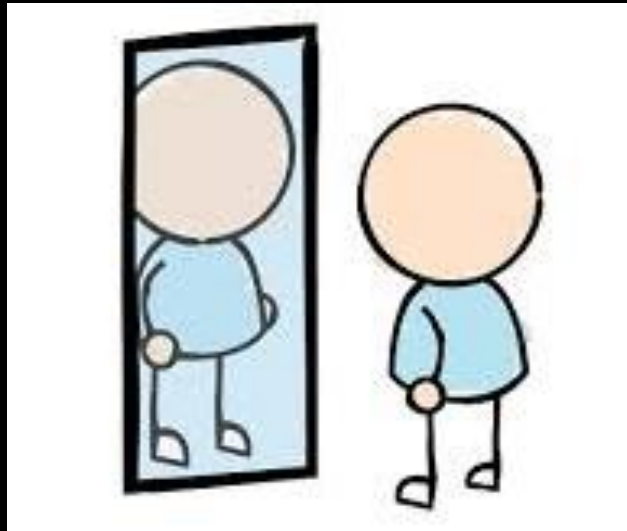
# Now it's your turn to drive. Start to finish.

**Your task**:

1. Pick your package and word-embedded vectors – it's between Gensim (Word2Vec) and SpaCy (GloVe).

2. Write code to **calculate the semantic similarity of two words** (e.g. *janky, ghetto*). "How similar in meaning?"

# 4 TAKE-AWAYS

1. **Google Ngram Viewer** is a quick 'n dirty tool for measuring word frequency change over time.

2. **Topic modeling** is a dimensionality reduction technique used to reveal "topics" in a document.

3. **Regular Expressions (RegEx) is the syntax you use** to do string matching, text cleaning, and token extraction.

4. **Word-embedded vectors** are decomposed matrices from a huge word matrix that tells you about word meaning.

# CHECK-IN:



1. So far, what is the most insightful thing you've learned during camp?
2. What is the one thing that's still the muddiest for you?