

Admin

- exams back tomorrow
- mini prj:
 - sign up to present
 - slides (PDF) due 2/23 9pm (nothing later!)
 - presses 2/24, 2/25
- semester project
 - proposal due 2/28 9pm (nothing later!)

Agenda

1. Sentiment Analysis
2. Web scraping
3. Python

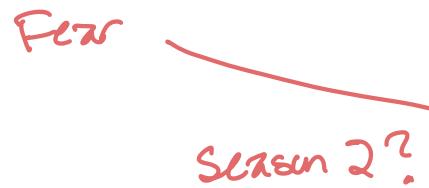
} data sources

1. Sentiment Analysis

↳ algorithm in Natural Language Processing (NLP)

Goal for this week:

- Should Laney keep watering the bear?
- ↳ do the vibes change over time?

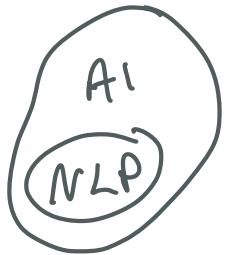


- figure art: ↗ tell us the mood of the episodes based on text descriptions

- ① ↗ - read ep description
- tell you if happy/sad

↗ good
↘ bad

NLP - ↗ spoken, written, humans with each other
⬇️ can't capture meaning
lacks precision, add nuance/sarcasm/jokes



↗ data science makes us better

meaning not important



Speech to text
one-to-one
translation

meaning matters

- care about the meaning behind what was said
(sentiment analysis)

Sentiment analysis

Ex) Polarity score -1 to +1
one score per phrase, sentence, etc.

nice weather!



+1 nice

$$y_2 = .5$$

fun mail :



+1 fun

$$y_2 = .5$$

SPRING BREAKER



$$y_2 = 0$$

Algorithm:-

- ignore caps, punc, stopwords
- add 1 for pos word
- sub 1 for neg word
- divide by # words

In Python: Text Blob

blob = TextBlob("nnnn")

blob.Sentiment.polarity -1 to +1

blob.Sentiment.subjectivity

0 to 1
fraction

What would make this algo better?

- emphasis words (very, really)

- slangy words (ex: fire, g2s, dope)
- different levels for good/bad
- text speak lol, rofl, lmao ;)
imho, imo, tbh
- negation words (not, don't)

2. web scraping

↳ assume: no CSV file of episodes
but, there is a wiki!

CSV file: structured data

```
~~, ~, ~~  
~~, ~, ~~  
~~, ~, ~~
```

HTML Crash Course

```
<html>  
  <title> My Page </title>  
  <body>  
    <h1> Header </h1>  
    <p> Hello! </p>  
    <img src='~~.jpg'>  
  </body>  
</html>
```

source of webpage

web page: semi-structured

- tags identify elements

`<~~>` `</~~>`
open close

- text in between tags
- created by humans
- Find the tags we need,
scrape the text in between



~~.html

The Bear Wiki

Season 1

ep1: 



ep2: 



- look in source for the tags in a description

<p class="description"> 
 </p>

- find all tags
- scrape text in between

3. Python

- sentiment analysis - how sophisticated?
- mini web scraping
- Scrape The Bear wiki, sentiment score for every episode of season 2.

Install:

- beautifulsoup4
- textblob
- wordcloud