

[Open in app ↗](#)

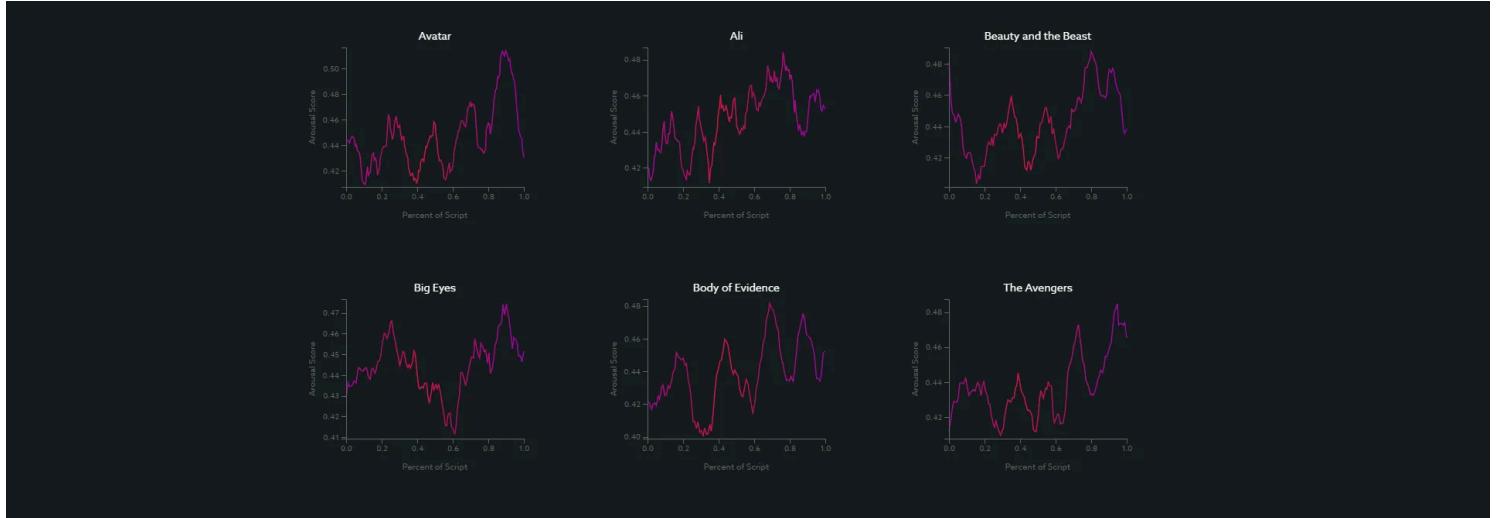
Search



Write



◆ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



Visualizing the Emotional Arcs of Movie Scripts Using Rule-Based Sentiment Analysis

How I used Python, D3 and Flask to create this [interactive visualization](#)



Nayomi Chibana · Follow

Published in Towards Data Science · 4 min read · Oct 17, 2019

206

3



...

Almost 72 years ago, acclaimed American writer Kurt Vonnegut came up with a novel method for [graphing the plot lines of stories](#) as part of his master's thesis in anthropology.

Although his work was ultimately rejected by the University of Chicago “because it was so simple and looked like too much fun,” according to Vonnegut, his overlooked contribution has received some renewed attention in the last few years after a group of researchers from the University of Vermont decided to use computational methods to test his hypothesis.

What they came up with were computer-generated story arcs for nearly 2,000 books in English, categorized into one of the six main storytelling arcs outlined in Vonnegut’s original thesis. These include “Rags to Riches” (rise), “Riches to Rags” (fall), “Man in a Hole” (fall then rise), “Icarus” (rise then fall), “Cinderella” (rise then fall then rise) and “Oedipus” (fall then rise then fall).

Their work deviated from Vonnegut’s in that they plotted the emotional trajectory of stories, not just their plot lines. To do so, they slid 10,000-word windows through each text to score the relative happiness of hundreds of points in the story, using a lexicon of 10,000 unique words scored on a nine-point scale of happiness, resulting in the hedonometer tool for sentiment analysis.

Using Arousal as a Proxy Measure for Action

Using a similar lexicon-based approach, I plotted the emotional arcs of more than a thousand movie scripts and used hierarchical clustering to group the most similar scripts.

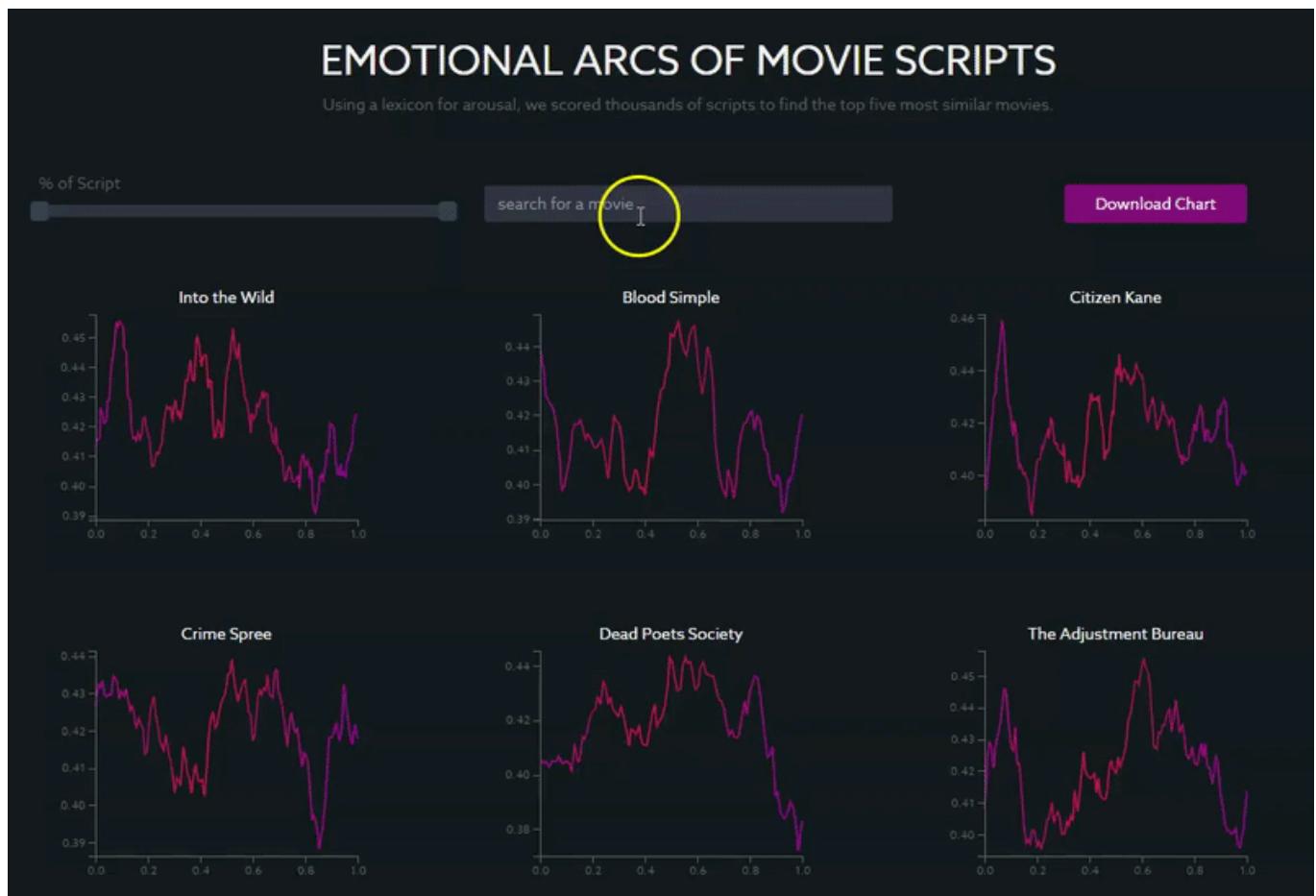
However, since previous research focuses mostly on positive and negative sentiment, rather than rising and falling action, the following method differs from UVM’s approach in that it uses the NRC Valence, Arousal and Dominance Lexicon, which has scores for 20,000 English words, as a proxy measure for action or conflict in a story. Specifically, I used the arousal

dimension to score words on a spectrum from “calm” or “passive” to “excited” or “active.”

The result is an interactive visualization that can be used to search for any script already published on the Internet Movie Database to visualize its emotional story arc (for arousal) and find five of the most similar movie scripts in terms of the story’s emotional trajectory, not content.

In a search for the blockbuster “Avatar,” for example, you can see that arousal peaks towards the end of the movie, at around 90% of the length of the script, which corresponds to the point in the movie with the most tension and conflict (the final confrontation between the Na’vi and the humans to determine the fate of the planet).

Other movies with similar structures, with a clear climactic end, include “Beauty and the Beast,” “The Avengers” (2012) and “Ali.”



Here is the process I followed to arrive at these results:

Scraping Texts and Scoring Using the VAD Lexicon

Using Beautiful Soup and a revised version of [Jeremy Kunn's script](#), I scraped all of the movies from the Internet Movie Database, which took about 20 minutes using the following script:

```
1 import os
2 from urllib.parse import quote
3
4 from bs4 import BeautifulSoup
5 import requests
6
7 BASE_URL = 'http://www.imsdb.com'
8 SCRIPTS_DIR = 'scripts'
9
10
11 def clean_script(text):
12     text = text.replace('Back to IMSDb', '')
13     text = text.replace(''''<b><!--',
14 </b>if (window!= top)
15 top.location.href=location.href
16 <b>// -->
17 </b>
18 '''', '')
19     text = text.replace(''''           Scanned by http://freemoviescripts.com
20                   Formatting by http://simplyscripts.home.att.net
21 ''', '')
22     return text.replace(r'\r', '')
23
24
25 def get_script(relative_link):
26     tail = relative_link.split('/')[-1]
27     print('fetching %s' % tail)
28     script_front_url = BASE_URL + quote(relative_link)
29     front_page_response = requests.get(script_front_url)
30     front_soup = BeautifulSoup(front_page_response.text, "html.parser")
31
32     try:
33         script_link = front_soup.find_all('p', align="center")[0].a['href']
34     except IndexError:
35         print('%s has no script :(' % tail)
36         return None, None
37
38     if script_link.endswith('.html'):
39         title = script_link.split('/')[-1].split(' Script')[0]
40         script_url = BASE_URL + script_link
41         script_soup = BeautifulSoup(requests.get(script_url).text, "html.parser")
42         script_text = script_soup.find_all('td', {'class': "scrtext"})[0].get_text()
43         script_text = clean_script(script_text)
44         return title, script_text
45     else:
```

```

46         print('%s is a pdf :(% % tail)
47         return None, None
48
49
50 if __name__ == "__main__":
51     response = requests.get('http://www.imsdb.com/all%20scripts/')
52     html = response.text
53
54     soup = BeautifulSoup(html, "html.parser")
55     paragraphs = soup.find_all('p')
56
57     for p in paragraphs:
58         relative_link = p.a['href']
59         title, script = get_script(relative_link)
60         if not script:
61             continue
62
63         with open(os.path.join(SCRIPTS_DIR, title.strip('.html') + '.txt'), 'w', encoding='utf-8') as outfile:
64             outfile.write(script)

```

download_all_scripts.py hosted with ❤ by GitHub

[view raw](#)

Then I created a new dictionary using the NRC VAD lexicon:

```

1 dict = pd.read_csv('./NRC-VAD-Lexicon-Aug2018Release/NRC-VAD-Lexicon-Aug2018Release/NRC-VAD-Lexicon.csv')
2 dict['Ranking'] = np.arange(1, len(dict)+1)
3 columnsTitles = ["Word", "Ranking", "Arousal", "Valence", "Dominance"]
4 dict = dict.reindex(columns=columnsTitles)
5 dict['Arousal'] = dict['Arousal'].astype(str)
6 newDict = dict.set_index('Word').T.to_dict('list')

```

arousal_lexicon.py hosted with ❤ by GitHub

[view raw](#)

And I adjusted the simple labMT usage script to calculate the arousal scores of a corpus, rather than happiness, by replacing the labMT dictionary and labMT vectors with my own:

```

1  def process():
2      windowSizes = [2000]
3      words = [x.lower() for x in re.findall(r"[\w\@\#\`\&]\*\-\/\[\=\;]+", raw_text_clean, f
4      lines = raw_text_clean.split("\n")
5      kwords = []
6      klines = []
7      for i in range(len(lines)):
8          if lines[i][0:3] != "<b>":
9              tmpwords = [x.lower() for x in re.findall(r"[\w\@\#\`\&]\*\-\/\[\=\;]+", lines[i], f
10             kwords.extend(tmpwords)
11             klines.extend([i for j in range(len(tmpwords))]))
12
13     for window in windowSizes:
14         breaks = [klines[window/10*i] for i in range(int(floor(float(len(klines))/window*10)))
15         breaks[0] = 0
16         f = open("word-vectors/" + str(window) + "/" + movie + "-breaks.csv", "w")
17         f.write(",".join(map(str, breaks)))
18         f.close()
19         chopper(kwords, labMT, labMTvector, "word-vectors/" + str(window) + "/" + movie + ".csv", minSize)
20
21         f = open("word-vectors/" + str(window) + "/" + movie + ".csv", "r")
22         fullVec = [list(map(int, line.split(","))) for line in f]
23         f.close()
24
25     # some movies are blank
26     if len(list(fullVec)) > 0:
27         if len(list(fullVec[0])) > 9:
28             precomputeTimeseries(fullVec, labMT, labMTvector, "timeseries/" + str(window) + "/" + movie)
29         else:
30             print("this movie is blank:")
31             print(movie.title)
32             movie.exclude = True
33             movie.excludeReason = "movie blank"

```

process.py hosted with ❤ by GitHub

[view raw](#)

Considering that movie scripts are shorter on average than books, I set the fixed window size at 1,000 words and slid it through each script to generate n arousal scores, or the number of points in the resulting time series.

Matrix Decomposition and Hierarchical Clustering

After implementing the revised version of simple labMT, which uses singular value decomposition under the hood to find a decomposition of stories onto an orthogonal basis of emotional arcs, I used linear interpolation to create word vectors of equal dimensions and used [scipy's hierarchical clustering](#) to find and group the most similar movie scripts based on the trajectory of their emotional arcs.

```
1 from scipy.cluster.hierarchy import fcluster
2
3 import scipy.cluster.hierarchy as hac
4
5 Z = hac.linkage(df_interpolate.iloc[:,0:266], method='ward', metric='euclidean')
6
7 # k Number of clusters I'd like to extract
8 results = fcluster(Z, k, criterion='maxclust')
```

cluster.py hosted with ❤ by GitHub

[view raw](#)

Using Ward's method and Euclidean distance as the distance metric, I was able to minimize the variance between clusters of movie scripts to arrive at the most accurate segmentation of observations.

Using Flask API to Score and D3 to Visualize Results

Finally, I created a Flask API to output arousal scores for any movie found in the Internet Movie Database and return the five most similar scripts.

Using a function to call the API and filter the data, I visualized the resulting output using D3's built-in d3.json() method for loading and parsing json data objects.

```

1  function updateData(){
2      d3.json('https://storyplotsapp.herokuapp.com/api', {
3          method:"POST",
4          body: JSON.stringify({
5              movie_title: ItemSelect
6          }),
7          headers: {"Content-type": "application/json; charset=UTF-8"}
8      }).then(function(data){
9          console.log(data);
10         // // Prepare and clean data
11         filteredData = {};
12         for (var movie in data) {
13             if (!data.hasOwnProperty(movie)) {
14                 continue;
15             }
16             filteredData[movie] = data[movie];
17             filteredData[movie].forEach(function(d){
18                 d["score"] = +d["score"];
19                 d["percent"] = +d["percent"];
20                 d["cluster3"] = +d["cluster3"];
21             });
22             console.log(filteredData);
23         }
24         d3.selectAll("svg").remove();
25
26         indexSearch = d3.keys(data).indexOf(ItemSelect);
27         var indexes = [0,1,2,3,4,5]
28         new_array = indexes.filter(function checkIndex(index) {
29             return index !== indexSearch;
30         });
31         console.log(new_array);
32
33         lineChart1 = new LineChart("#chart-area1", d3.keys(data)[indexSearch]);
34         lineChart2 = new LineChart("#chart-area2", d3.keys(data)[new_array[0]]);
35         lineChart3 = new LineChart("#chart-area3", d3.keys(data)[new_array[1]]);
36         lineChart4 = new LineChart("#chart-area4", d3.keys(data)[new_array[2]]);
37         lineChart5 = new LineChart("#chart-area5", d3.keys(data)[new_array[3]]);
38         lineChart6 = new LineChart("#chart-area6", d3.keys(data)[new_array[4]]);
39     })
40 }

```

updatedata.js hosted with ❤ by GitHub

[view raw](#)

You can play around with the final interactive visualization [here](#) and view all the notebooks and code [here](#). Let me know what you think below!

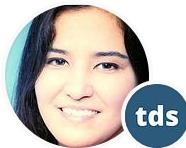
Machine Learning

Data Visualization

Data Science

Sentiment Analysis

Storytelling



Written by Nayomi Chibana

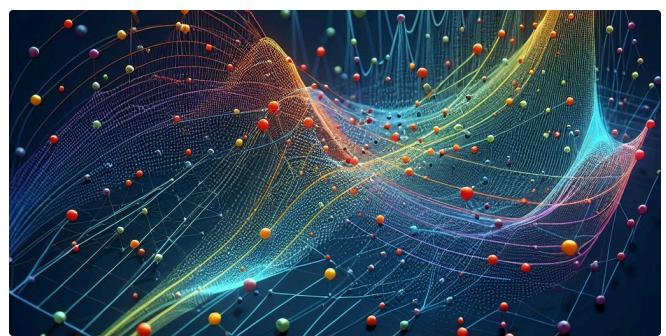
30 Followers · Writer for Towards Data Science

Follow



Data journalist interested in all things at the intersection of news, data and tech.

More from Nayomi Chibana and Towards Data Science





Nayomi Chibana in Towards Data Science



Tim Sumner in Towards Data Science

How Do Psychometric Test Results Vary Across Age, Race and...

Most of us have taken a personality quiz at one point or another. Whether it was for a job...

5 min read · Jun 28, 2019



154



A New Coefficient of Correlation

What if you were told there exists a new way to measure the relationship between two...

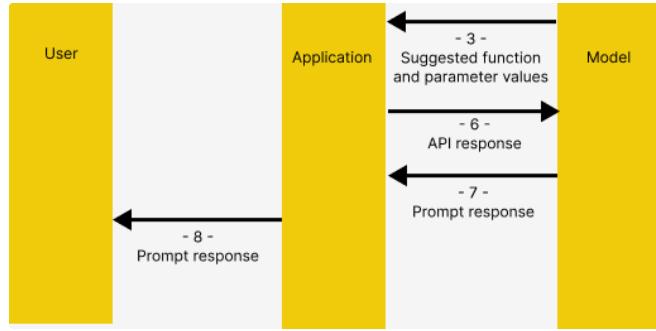
10 min read · Mar 31, 2024



3.5K



44



Youness Mansar in Towards Data Science

Meet the NiceGUI: Your Soon-to-be Favorite Python UI Library

Build custom web apps easily and quickly

8 min read · Apr 16, 2024



1.2K



Julian Yip in Towards Data Science

Build Autonomous AI Agents with Function Calling

Transform your chatbot into an agent that can interact with external APIs

11 min read · Apr 2, 2024



1.1K



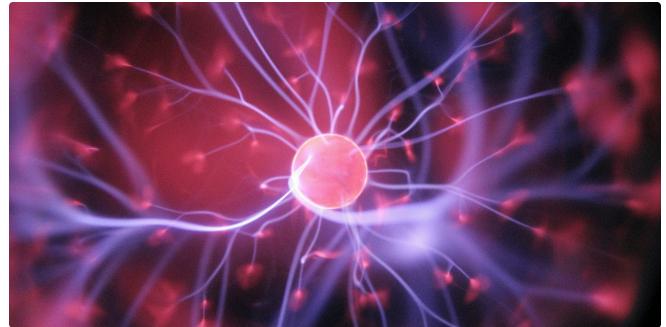
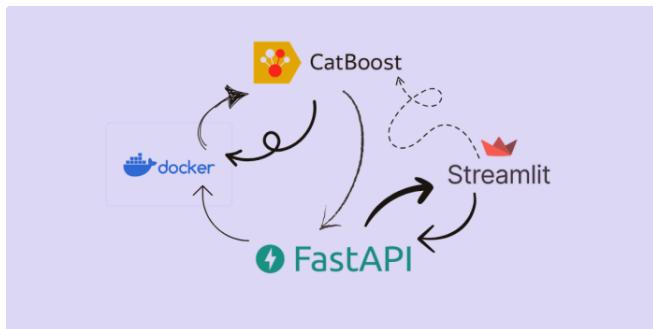
8



[See all from Nayomi Chibana](#)

[See all from Towards Data Science](#)

Recommended from Medium



 Ramazan Olmez

End-to-End Machine Learning Project: Churn Prediction

The main objective of this article is to develop an end-to-end machine learning project. For...

18 min read · Feb 22, 2024

 538  4

 + 

 Rosaria Silipo  in Low Code for Data Science

Is Data Science dead?

In the last six months I have heard this question thousands of time: “Is data science...

6 min read · Mar 11, 2024

 2.1K  47

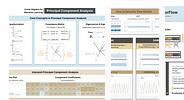
 + 

Lists



Predictive Modeling w/ Python

20 stories · 1161 saves



Practical Guides to Machine Learning

10 stories · 1402 saves



Natural Language Processing

1434 stories · 929 saves



data science and AI

40 stories · 149 saves



 Bale Chen in Towards Data Science

Emojis Aid Social Media Sentiment Analysis: Stop Cleaning Them Out!

Leverage emojis in social media sentiment analysis to improve accuracy.

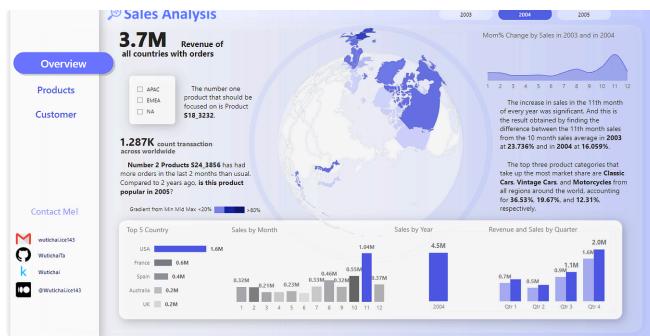
14 min read · Jan 31, 2023

 332

 2



...



 Wutichai Ta

My First Data Analyst Project (Portfolio)—Python: Sales Analysis

This is my first data analysis project using Python. It was the beginning of my journey...

5 min read · Apr 23, 2024

 275

 8



...

 Ubaid Shah

Twitter Sentiment Analysis — A Step by Guide

In today's digital age, social media platforms like Twitter have become a goldmine for...

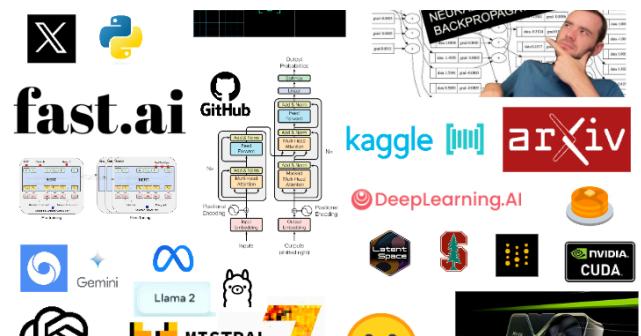
6 min read · Nov 8, 2023

 31





...



 Benedict Neo in bitgrit Data Science Publication

Roadmap to Learn AI in 2024

A free curriculum for hackers and programmers to learn AI

11 min read · Mar 11, 2024

 10.9K

 119



...

See more recommendations