

## **Introduction**

### **Sentiment Analysis**

“Sentiment Analysis is the interpretation and classification of emotions (positive, negative, or neutral) within data using text analysis techniques. It allows business to identify customer sentiment towards products or services in online conversations and feedback.” (Tran, et al. 2022).

The entertainment and multimedia sector occupies a large share of the market today, and an important big part of this sector goes to movie production.

There are a group of websites that offer platforms dedicated to review the movies, such as IMDB and Rotten Tomatoes, and these websites also provide a list of preferences for films in addition to many other services.

### **The task**

In this research, I will do simple partial sentiment analysis to evaluate a sample of movies that share the same genre using the positive and negative reviews dataset of Rotten Tomatoes. The study could give a glance at the reaction of the population towards a specific type of film and that could be used to achieve two main goals: the first goal is to nominate some other movie names from the same genre, and the other goal is to help the people in the movie production industry to understand the reasons of success or failure of a movie and encourage them to focus on the formula of the success. I applied the unsupervised learning method (word2vec) to analyze 3 movies from the (organized crime) movies genre and those are (The Godfather, The Goodfellas, and Scarface).

### **Dataset and preprocessing**

I chose to study the comments of the reviewers from the Rotten Tomato website, and I used the data set provided at Github.com (<https://github.com/nicolas-gervais/rotten-tomatoes-dataset>) published by the author Nicolas Gervais.

I started with loading the number of required packages such as nltk, genism, and panda, then added a function to recognize the punctuation in the text, after that used the panda package to call, read and process the dataset file. I examined the file reading function by obtaining the first and last 10 records of the dataset file.

After that, I converted the text to lower case so I can get rid of all the capitalized or upper case text, and get the source text unified. Then did the tokenization to split the text into pieces of words (Unigram), ignore too short or too long tokens, and create unique words dictionary.

## **The method**

The next function was the creation of an embedding model from text using word2vec method in genism. The word2vec method algorithm is a technique used to process the natural languages, and it can learn the relationships of the words of a large text source, then generate synonyms and vocab suggestions or solutions as needed.

Then tried to create a bag of words to measure the times a specific word was repeated, but Jupyter was either collapsing or generating the message of (IOPub data rate exceeded) (**Fig.08, Appendix 01**), after a long time of trying to solve this issue, I decided to go with another solution which is writing a function to search within the original file and get the times of a specific word occurrence (in both upper and lower cases) and this function worked perfectly. After that, I tried to explore the text and find the similarity between the negative and positive words.

I searched and brought some of the most common negative and positive adjectives and added them to a function that measures the similarity of a movie name alongside the adjective and provides a value of similarity. Here is a sample of the list of adjectives that I used within the function ('weak', 'wacky', 'stupid', 'silly', 'average', 'brilliant', 'charismatic', 'charming', 'clever' ).

First, I used the words similarity values for each of the selected movies by calculating the sum of the values of the positive adjective then extracting the average value, and did the same with the negative values.

Second, I extracted the minimum and maximum of each of the positive and negative values. Third I tried to find the times that each of the movie names mentioned in the dataset. Finally, I generated some visualizations for the results.

## **Findings**

The similarity of the adjectives towards each of the movies, after calculating the average of the positive and negative similarity in addition to the calculation of the minimum and maximum values, and the movie name repetition times within the reviewer comments are as follows:

### The Godfather

	Average	Min	Max
Positive	0.111024035	-0.014134126	0.47434053
Negative	0.01386544	-0.089685224	0.11394882

Times of Repetition 142

### The Goodfellas

	Average	Min	Max
Positive	0.125275256	-0.000528622	0.34355572
Negative	0.098974584	-0.033701137	0.20421866

Times of Repetition 70

### Scarface

	Average	Min	Max
Positive	0.117413777	-0.004182309	0.36156192
Negative	0.0722755	-0.008026212	0.23615609

Times of Repetition 47

## Results

It is noticeable from the up-mentioned result that the chosen movies achieved a high value of similarity regarding positive reviews compared to the negative reviews. The movie name repetition in the review comments shows that The Godfather movie was mentioned (147 times), the Goodfellas (70 times) while Scarface movie achieved the least value of mentioning (47 times).

The results show that the godfather movie could be used as a major guide or measuring reference for movie recommendations from the same genre so whoever shows some positive review or reaction

towards the Godfather movie could be provided with a list of movies from the same classification such as The Goodfellas, etc.

The results also reflect the success of the Godfather movie formula, which could be followed as a guide in the production of any future movies.

This was a quick study based on limited knowledge of the tools, however, it could be applied to wider genres or classifications in the future and get more accurate to build a robust recommendation system based on sentiment analysis.

## Reference:

Tran, D.D., Nguyen, T.T., Dao, T. (2022). *Sentiment Analysis of Movie Reviews Using Machine Learning Techniques. Proceedings of Sixth International Congress on Information and Communication Technology*. Lecture Notes in Networks and Systems, vol 235. Springer, Singapore. [https://doi.org/10.1007/978-981-16-2377-6\\_34](https://doi.org/10.1007/978-981-16-2377-6_34)

## Appendix 01

### Create a word embedding model and explore patterns in the source dataset text- with PANDAS

```
In [1]: # Import packages
```

```
In [2]: import nltk
import re
import gensim
import csv
import nltk.tokenize
import pandas as pd
nltk.download('punkt', quiet=True)
from sklearn.manifold import TSNE
import numpy as np
import pandas as pd
```

```
In [3]: #FUNCTION TO RECOGNISE PUNCTUATION
PUNCT_RE = re.compile(r'^\w\s]+$')
def is_punct(string):
    return PUNCT_RE.match(string) is not None
```

```
In [4]: df = pd.read_csv("F:\RottenTomato.csv", encoding='utf-8')
```

```
In [5]: #Returns the first N rows
df.head(10)
```

Out[5]:

	Freshness	Review
0	fresh	Manakamana doesn't answer any questions, yet ...
1	fresh	Wiffully offensive and powered by a chest-thu...
2	rotten	It would be difficult to imagine material mor...
3	rotten	Despite the gusto its star brings to the role...
4	rotten	If there was a good idea at the core of this ...
5	rotten	Gleeson goes the Hallmark Channel route, dama...
6	fresh	It was the height of satire in 1976: dark as ...
7	rotten	Everyone in "The Comedian" deserves a better ...
8	rotten	Actor encourages grumpy Christians to embrace...
9	fresh	Slight, contained, but ineffably soulful.

Fig.01

```
In [6]: #Returns the Last N rows
df.tail(10)
```

Out[6]:

	Freshness	Review
479990	fresh	Ryan Coogler's sequel follows a "Rocky" road
479991	fresh	Disney brilliantly executed another film in t...
479992	rotten	Director John Crowley overestimates the comed...
479993	rotten	Here's a sobering thought: If every war gets ...
479994	rotten	Roland Joffe's deeply ridiculous movie is cau...
479995	rotten	Zemeckis seems unable to admit that the motio...
479996	fresh	Movies like The Kids Are All Right -- beautif...
479997	rotten	Film-savvy audiences soon will catch onto Win...
479998	fresh	An odd yet enjoyable film.
479999	fresh	No other animation studio, even our beloved P...

Fig.02

```
In [7]: #CONVERT STRINGS TO LOWERCASE
```

```
reviews = df['Review'].str.lower()
```

```
In [8]: reviews
```

```
Out[8]: 0      manakamana doesn't answer any questions, yet ...
        1      wilfully offensive and powered by a chest-thu...
        2      it would be difficult to imagine material mor...
        3      despite the gusto its star brings to the role...
        4      if there was a good idea at the core of this ...
        ...
        479995     zemeckis seems unable to admit that the motio...
        479996     movies like the kids are all right -- beautif...
        479997     film-savvy audiences soon will catch onto win...
        479998                an odd yet enjoyable film.
        479999     no other animation studio, even our beloved p...
        Name: Review, Length: 480000, dtype: object
```

Fig.03

```
In [9]: #LIST THE VALUES IN COLUMN
```

```
reviews = df['Review'].str.lower().values
reviews
```

```
Out[9]: array([" manakamana doesn't answer any questions, yet makes its point: nepal, like the rest of our planet, is a picturesque but
far from peaceable kingdom.",
      " wilfully offensive and powered by a chest-thumping machismo, but it's good clean fun.",
      ' it would be difficult to imagine material more wrong for spade than lost & found.',
      ...,
      " film-savvy audiences soon will catch onto winterbottom's attempts to earn neo-noir credentials as patel and apte's cha
racters go on the run. patel overdoes his character's pragmatic stoicism, and i eventually tired of the waiting game...",
      ' an odd yet enjoyable film.',
      ' no other animation studio, even our beloved pixar, can quite replicate the magic, wonder and weirdness [of studio ghib
l].'],
      dtype=object)
```

Fig.04

```
In [11]: #Tokenize text ()
reviews = [nltk.word_tokenize(review) for review in reviews]
```

```
In [12]: reviews
```

```
Out[12]: [['manakamana',
          'does',
          "n't",
          'answer',
          'any',
          'questions',
          ',',
          'yet',
          'makes',
          'its',
          'point',
          ':',
          'nepal',
          ',',
          'like',
          'the',
          'rest',
          'of',
          'our',
          '...']]
```

Fig.05

```
In [13]: #Import gensim's packages
```

```
from gensim.utils import simple_preprocess
from gensim import corpora
#convert document into list of lowercase tokens, ignore too short or too long tokens, create unique words dictionary
dictionary = corpora.Dictionary(simple_preprocess(line, deacc = True) for line in open('F:\RottenTomato.csv'))
print (dictionary.token2id)
```

```
{'freshness': 0, 'review': 1, 'answer': 2, 'any': 3, 'but': 4, 'doesn': 5, 'far': 6, 'fresh': 7, 'from': 8, 'is': 9, 'its': 10, 'kingdom': 11, 'like': 12, 'makes': 13, 'manakamana': 14, 'nepal': 15, 'of': 16, 'our': 17, 'peaceable': 18, 'picturesque': 19, 'planet': 20, 'point': 21, 'questions': 22, 'rest': 23, 'the': 24, 'yet': 25, 'and': 26, 'by': 27, 'chest': 28, 'clean': 29, 'fun': 30, 'good': 31, 'it': 32, 'machismo': 33, 'offensive': 34, 'powered': 35, 'thumping': 36, 'wilfully': 37, 'better': 38, 'difficult': 39, 'for': 40, 'found': 41, 'imagine': 42, 'lost': 43, 'material': 44, 'more': 45, 'rotten': 46, 'spade': 47, 'than': 48, 'to': 49, 'would': 50, 'wrong': 51, 'brings': 52, 'despite': 53, 'discovery': 54, 'gusto': 55, 'hard': 56, 'hector': 57, 'on': 58, 'ride': 59, 'role': 60, 'shotgun': 61, 'star': 62, 'voyage': 63, 'an': 64, 'arson': 65, 'at': 66, 'bad': 67, 'been': 68, 'buried': 69, 'core': 70, 'dog': 71, 'film': 72, 'flatulence': 73, 'idea': 74, 'if': 75, 'in': 76, 'jokes': 77, 'pile': 78, 'plot': 79, 'puns': 80, 'related': 81, 'ridiculous': 82, 'serial': 83, 'there': 84, 'this': 85, 'unusually': 86, 'was': 87, 'channel': 88, 'curious': 89, 'damaging': 90, 'entry': 91, 'gleeson': 92, 'goes': 93, 'hallmark': 94, 'intermittently': 95, 'route': 96, 'subgenre': 97, 'time': 98, 'travel': 99, 'absurd': 100, 'as': 101, 'close': 102, 'dark': 103, 'era': 104, 'height': 105, 'hell': 106, 'jerry': 107, 'nowhere': 108, 'objective': 109, 'patently': 110, 'reality': 111, 'satire': 112, 'somewhere': 113, 'springer': 114, 'surely': 115, 'surpassed': 116, 'better': 117, 'comedian': 118, 'deserves': 119, 'everyone': 120, 'movie': 121, 'actor': 122, 'christians': 123, 'embrace': 124, 'encourages': 125, 'grumpy': 126, 'season': 127, 'contained': 128, 'ineffably': 129, 'slight': 130, 'soulful': 131, 'addresses': 132, 'angle': 133, 'approaching': 134, 'bell': 135, 'both': 136, 'conceivable': 137, 'debut': 138, 'enhancement': 139, 'every': 140, 'feature': 141, 'humor': 142, 'intelligence': 143, 'issue': 144, 'performance': 145, 'subject': 146, 'with': 147, 'achieves': 148, 'almost': 149, 'considerable': 150, 'director': 151, 'entirely': 152, 'except': 153, 'halloween': 154, 'irresistibly': 155, 'means': 156, 'music': 157, 'power': 158, 'score': 159, 'shamelessly': 160, 'through': 161, 'visual': 162, 'zingy': 163, 'action': 164, '...': 165, '...': 166, '...': 167, '...': 168, '...': 169, '...': 170, '...': 171, '...': 172, '...': 173, '...': 174, '...': 175, '...': 176, '...': 177, '...': 178, '...': 179, '...': 180, '...': 181, '...': 182, '...': 183, '...': 184, '...': 185, '...': 186, '...': 187, '...': 188, '...': 189, '...': 190, '...': 191, '...': 192, '...': 193, '...': 194, '...': 195, '...': 196, '...': 197, '...': 198, '...': 199, '...': 200, '...': 201, '...': 202, '...': 203, '...': 204, '...': 205, '...': 206, '...': 207, '...': 208, '...': 209, '...': 210, '...': 211, '...': 212, '...': 213, '...': 214, '...': 215, '...': 216, '...': 217, '...': 218, '...': 219, '...': 220, '...': 221, '...': 222, '...': 223, '...': 224, '...': 225, '...': 226, '...': 227, '...': 228, '...': 229, '...': 230, '...': 231, '...': 232, '...': 233, '...': 234, '...': 235, '...': 236, '...': 237, '...': 238, '...': 239, '...': 240, '...': 241, '...': 242, '...': 243, '...': 244, '...': 245, '...': 246, '...': 247, '...': 248, '...': 249, '...': 250, '...': 251, '...': 252, '...': 253, '...': 254, '...': 255, '...': 256, '...': 257, '...': 258, '...': 259, '...': 260, '...': 261, '...': 262, '...': 263, '...': 264, '...': 265, '...': 266, '...': 267, '...': 268, '...': 269, '...': 270, '...': 271, '...': 272, '...': 273, '...': 274, '...': 275, '...': 276, '...': 277, '...': 278, '...': 279, '...': 280, '...': 281, '...': 282, '...': 283, '...': 284, '...': 285, '...': 286, '...': 287, '...': 288, '...': 289, '...': 290, '...': 291, '...': 292, '...': 293, '...': 294, '...': 295, '...': 296, '...': 297, '...': 298, '...': 299, '...': 300, '...': 301, '...': 302, '...': 303, '...': 304, '...': 305, '...': 306, '...': 307, '...': 308, '...': 309, '...': 310, '...': 311, '...': 312, '...': 313, '...': 314, '...': 315, '...': 316, '...': 317, '...': 318, '...': 319, '...': 320, '...': 321, '...': 322, '...': 323, '...': 324, '...': 325, '...': 326, '...': 327, '...': 328, '...': 329, '...': 330, '...': 331, '...': 332, '...': 333, '...': 334, '...': 335, '...': 336, '...': 337, '...': 338, '...': 339, '...': 340, '...': 341, '...': 342, '...': 343, '...': 344, '...': 345, '...': 346, '...': 347, '...': 348, '...': 349, '...': 350, '...': 351, '...': 352, '...': 353, '...': 354, '...': 355, '...': 356, '...': 357, '...': 358, '...': 359, '...': 360, '...': 361, '...': 362, '...': 363, '...': 364, '...': 365, '...': 366, '...': 367, '...': 368, '...': 369, '...': 370, '...': 371, '...': 372, '...': 373, '...': 374, '...': 375, '...': 376, '...': 377, '...': 378, '...': 379, '...': 380, '...': 381, '...': 382, '...': 383, '...': 384, '...': 385, '...': 386, '...': 387, '...': 388, '...': 389, '...': 390, '...': 391, '...': 392, '...': 393, '...': 394, '...': 395, '...': 396, '...': 397, '...': 398, '...': 399, '...': 400, '...': 401, '...': 402, '...': 403, '...': 404, '...': 405, '...': 406, '...': 407, '...': 408, '...': 409, '...': 410, '...': 411, '...': 412, '...': 413, '...': 414, '...': 415, '...': 416, '...': 417, '...': 418, '...': 419, '...': 420, '...': 421, '...': 422, '...': 423, '...': 424, '...': 425, '...': 426, '...': 427, '...': 428, '...': 429, '...': 430, '...': 431, '...': 432, '...': 433, '...': 434, '...': 435, '...': 436, '...': 437, '...': 438, '...': 439, '...': 440, '...': 441, '...': 442, '...': 443, '...': 444, '...': 445, '...': 446, '...': 447, '...': 448, '...': 449, '...': 450, '...': 451, '...': 452, '...': 453, '...': 454, '...': 455, '...': 456, '...': 457, '...': 458, '...': 459, '...': 460, '...': 461, '...': 462, '...': 463, '...': 464, '...': 465, '...': 466, '...': 467, '...': 468, '...': 469, '...': 470, '...': 471, '...': 472, '...': 473, '...': 474, '...': 475, '...': 476, '...': 477, '...': 478, '...': 479, '...': 480, '...': 481, '...': 482, '...': 483, '...': 484, '...': 485, '...': 486, '...': 487, '...': 488, '...': 489, '...': 490, '...': 491, '...': 492, '...': 493, '...': 494, '...': 495, '...': 496, '...': 497, '...': 498, '...': 499, '...': 500, '...': 501, '...': 502, '...': 503, '...': 504, '...': 505, '...': 506, '...': 507, '...': 508, '...': 509, '...': 510, '...': 511, '...': 512, '...': 513, '...': 514, '...': 515, '...': 516, '...': 517, '...': 518, '...': 519, '...': 520, '...': 521, '...': 522, '...': 523, '...': 524, '...': 525, '...': 526, '...': 527, '...': 528, '...': 529, '...': 530, '...': 531, '...': 532, '...': 533, '...': 534, '...': 535, '...': 536, '...': 537, '...': 538, '...': 539, '...': 540, '...': 541, '...': 542, '...': 543, '...': 544, '...': 545, '...': 546, '...': 547, '...': 548, '...': 549, '...': 550, '...': 551, '...': 552, '...': 553, '...': 554, '...': 555, '...': 556, '...': 557, '...': 558, '...': 559, '...': 560, '...': 561, '...': 562, '...': 563, '...': 564, '...': 565, '...': 566, '...': 567, '...': 568, '...': 569, '...': 570, '...': 571, '...': 572, '...': 573, '...': 574, '...': 575, '...': 576, '...': 577, '...': 578, '...': 579, '...': 580, '...': 581, '...': 582, '...': 583, '...': 584, '...': 585, '...': 586, '...': 587, '...': 588, '...': 589, '...': 590, '...': 591, '...': 592, '...': 593, '...': 594, '...': 595, '...': 596, '...': 597, '...': 598, '...': 599, '...': 600, '...': 601, '...': 602, '...': 603, '...': 604, '...': 605, '...': 606, '...': 607, '...': 608, '...': 609, '...': 610, '...': 611, '...': 612, '...': 613, '...': 614, '...': 615, '...': 616, '...': 617, '...': 618, '...': 619, '...': 620, '...': 621, '...': 622, '...': 623, '...': 624, '...': 625, '...': 626, '...': 627, '...': 628, '...': 629, '...': 630, '...': 631, '...': 632, '...': 633, '...': 634, '...': 635, '...': 636, '...': 637, '...': 638, '...': 639, '...': 640, '...': 641, '...': 642, '...': 643, '...': 644, '...': 645, '...': 646, '...': 647, '...': 648, '...': 649, '...': 650, '...': 651, '...': 652, '...': 653, '...': 654, '...': 655, '...': 656, '...': 657, '...': 658, '...': 659, '...': 660, '...': 661, '...': 662, '...': 663, '...': 664, '...': 665, '...': 666, '...': 667, '...': 668, '...': 669, '...': 670, '...': 671, '...': 672, '...': 673, '...': 674, '...': 675, '...': 676, '...': 677, '...': 678, '...': 679, '...': 680, '...': 681, '...': 682, '...': 683, '...': 684, '...': 685, '...': 686, '...': 687, '...': 688, '...': 689, '...': 690, '...': 691, '...': 692, '...': 693, '...': 694, '...': 695, '...': 696, '...': 697, '...': 698, '...': 699, '...': 700, '...': 701, '...': 702, '...': 703, '...': 704, '...': 705, '...': 706, '...': 707, '...': 708, '...': 709, '...': 710, '...': 711, '...': 712, '...': 713, '...': 714, '...': 715, '...': 716, '...': 717, '...': 718, '...': 719, '...': 720, '...': 721, '...': 722, '...': 723, '...': 724, '...': 725, '...': 726, '...': 727, '...': 728, '...': 729, '...': 730, '...': 731, '...': 732, '...': 733, '...': 734, '...': 735, '...': 736, '...': 737, '...': 738, '...': 739, '...': 740, '...': 741, '...': 742, '...': 743, '...': 744, '...': 745, '...': 746, '...': 747, '...': 748, '...': 749, '...': 750, '...': 751, '...': 752, '...': 753, '...': 754, '...': 755, '...': 756, '...': 757, '...': 758, '...': 759, '...': 760, '...': 761, '...': 762, '...': 763, '...': 764, '...': 765, '...': 766, '...': 767, '...': 768, '...': 769, '...': 770, '...': 771, '...': 772, '...': 773, '...': 774, '...': 775, '...': 776, '...': 777, '...': 778, '...': 779, '...': 780, '...': 781, '...': 782, '...': 783, '...': 784, '...': 785, '...': 786, '...': 787, '...': 788, '...': 789, '...': 790, '...': 791, '...': 792, '...': 793, '...': 794, '...': 795, '...': 796, '...': 797, '...': 798, '...': 799, '...': 800, '...': 801, '...': 802, '...': 803, '...': 804, '...': 805, '...': 806, '...': 807, '...': 808, '...': 809, '...': 810, '...': 811, '...': 812, '...': 813, '...': 814, '...': 815, '...': 816, '...': 817, '...': 818, '...': 819, '...': 820, '...': 821, '...': 822, '...': 823, '...': 824, '...': 825, '...': 826, '...': 827, '...': 828, '...': 829, '...': 830, '...': 831, '...': 832, '...': 833, '...': 834, '...': 835, '...': 836, '...': 837, '...': 838, '...': 839, '...': 840, '...': 841, '...': 842, '...': 843, '...': 844, '...': 845, '...': 846, '...': 847, '...': 848, '...': 849, '...': 850, '...': 851, '...': 852, '...': 853, '...': 854, '...': 855, '...': 856, '...': 857, '...': 858, '...': 859, '...': 860, '...': 861, '...': 862, '...': 863, '...': 864, '...': 865, '...': 866, '...': 867, '...': 868, '...': 869, '...': 870, '...': 871, '...': 872, '...': 873, '...': 874, '...': 875, '...': 876, '...': 877, '...': 878, '...': 879, '...': 880, '...': 881, '...': 882, '...': 883, '...': 884, '...': 885, '...': 886, '...': 887, '...': 888, '...': 889, '...': 890, '...': 891, '...': 892, '...': 893, '...': 894, '...': 895, '...': 896, '...': 897, '...': 898, '...': 899, '...': 900, '...': 901, '...': 902, '...': 903, '...': 904, '...': 905, '...': 906, '...': 907, '...': 908, '...': 909, '...': 910, '...': 911, '...': 912, '...': 913, '...': 914, '...': 915, '...': 916, '...': 917, '...': 918, '...': 919, '...': 920, '...': 921, '...': 922, '...': 923, '...': 924, '...': 925, '...': 926, '...': 927, '...': 928, '...': 929, '...': 930, '...': 931, '...': 932, '...': 933, '...': 934, '...': 935, '...': 936, '...': 937, '...': 938, '...': 939, '...': 940, '...': 941, '...': 942, '...': 943, '...': 944, '...': 945, '...': 946, '...': 947, '...': 948, '...': 949, '...': 950, '...': 951, '...': 952, '...': 953, '...': 954, '...': 955, '...': 956, '...': 957, '...': 958, '...': 959, '...': 960, '...': 961, '...': 962, '...': 963, '...': 964, '...': 965, '...': 966, '...': 967, '...': 968, '...': 969, '...': 970, '...': 971, '...': 972, '...': 973, '...': 974, '...': 975, '...': 976, '...': 977, '...': 978, '...': 979, '...': 980, '...': 981, '...': 982, '...': 983, '...': 984, '...': 985, '...': 986, '...': 987, '...': 988, '...': 989, '...': 990, '...': 991, '...': 992, '...': 993, '...': 994, '...': 995, '...': 996, '...': 997, '...': 998, '...': 999, '...': 1000, '...': 1001, '...': 1002, '...': 1003, '...': 1004, '...': 1005, '...': 1006, '...': 1007, '...': 1008, '...': 1009, '...': 1010, '...': 1011, '...': 1012, '...': 1013, '...': 1014, '...': 1015, '...': 1016, '...': 1017, '...': 1018, '...': 1019, '...': 1020, '...': 1021, '...': 1022, '...': 1023, '...': 1024, '...': 1025, '...': 1026, '...': 1027, '...': 1028, '...': 1029, '...': 1030, '...': 1031, '...': 1032, '...': 1033, '...': 1034, '...': 1035, '...': 1036, '...': 1037, '...': 1038, '...': 1039, '...': 1040, '...': 1041, '...': 1042, '...': 1043, '...': 1044, '...': 1045, '...': 1046, '...': 1047, '...': 1048, '...': 1049, '...': 1050, '...': 1051, '...': 1052, '...': 1053, '...': 1054, '...': 1055, '...': 1056, '...': 1057, '...': 1058, '...': 1059, '...': 1060, '...': 1061, '...': 1062, '...': 1063, '...': 1064, '...': 1065, '...': 1066, '...': 1067, '...': 1068, '...': 1069, '...': 1070, '...': 1071, '...': 1072, '...': 1073, '...': 1074, '...': 1075, '...': 1076, '...': 1077, '...': 1078, '...': 1079, '...': 1080, '...': 1081, '...': 1082, '...': 1083, '...': 1084, '...': 1085, '...': 1086, '...': 1087, '...': 1088, '...': 1089, '...': 1090, '...': 1091, '...': 1092, '...': 1093, '...': 1094, '...': 1095, '...': 1096, '...': 1097, '...': 1098, '...': 1099, '...': 1100, '...': 1101, '...': 1102, '...': 1103, '...': 1104, '...': 1105, '...': 1106, '...': 1107, '...': 1108, '...': 1109, '...': 1110, '...': 1111, '...': 1112, '...': 1113, '...': 1114, '...': 1115, '...': 1116, '...': 1117, '...': 1118, '...': 1119, '...': 1120, '...': 1121, '...': 1122, '...': 1123, '...': 1124, '...': 1125, '...': 1126, '...': 1127, '...': 1128, '...': 1129, '...': 1130, '...': 1131, '...': 1132, '...': 1133, '...': 1134, '...': 1135, '...': 1136, '...': 1137, '...': 1138, '...': 1139, '...': 1140, '...': 1141, '...': 1142, '...': 1143, '...': 1144, '...': 1145, '...': 1146, '...': 1147, '...': 1148, '...': 1149, '...': 1150, '...': 1151, '...': 1152, '...': 1153, '...': 1154, '...': 1155, '...': 1156, '...': 1157, '...': 1158, '...': 1159, '...': 1160, '...': 1161, '...': 1162, '...': 1163, '...': 1164, '...': 1165, '...': 1166, '...': 1167, '...': 1168, '...': 1169, '...': 1170, '...': 1171, '...': 1172, '...': 1173, '...': 1174, '...': 1175, '...': 1176, '...': 1177, '...': 1178, '...': 1179, '...': 1180, '...': 1181, '...': 1182, '...': 1183, '...': 1184, '...': 1185, '...': 1186, '...': 1187, '...': 1188, '...': 1189, '...': 1190, '...': 1191, '...': 1192, '...': 1193, '...': 1194, '...': 1195, '...': 1196, '...': 1197, '...': 1198, '...': 1199, '...': 1200, '...': 1201, '...': 1202, '...': 1203, '...': 1204, '...': 1205, '...': 1206, '...': 1207, '...': 1208, '...': 1209, '...': 1210, '...': 1211, '...': 1212, '...': 1213, '...': 1214, '...': 1215, '...': 1216, '...': 1217, '...': 1218, '...': 1219, '...': 1220, '...': 1221, '...': 1222, '...': 1223, '...': 1224, '...': 1225, '...': 1226, '...': 1227, '...': 1228, '...': 1229, '...': 1230, '...': 1231, '...': 1232, '...': 1233, '...': 1234, '...': 1235, '...': 1236, '...': 1237, '...': 1238, '...': 1239, '...': 1240, '...': 1241, '...': 1242, '...': 1243, '...': 1244, '...': 1245, '...': 1246, '...': 1247, '...': 1248, '...': 1249, '...': 1250, '...': 1251, '...': 1252, '...': 1253, '...': 1254, '...': 1255, '...': 1256, '...': 1257, '...': 1258, '...': 1259, '...': 1260, '...': 1261, '...': 1262, '...': 1263, '...': 1264, '...': 1265, '...': 1266, '...': 1267, '...': 1268, '...': 1269, '...': 1270, '...': 1271, '...': 1272, '...': 1273, '...': 1274, '...': 1275, '...': 1276, '...': 1277, '...': 1278, '...': 1279, '...': 1280, '...': 1281, '...': 1282, '...': 1283, '...': 1284, '...': 1285, '...': 1286, '...': 1287, '...': 1288, '...': 1289, '...': 1290, '...': 1291, '...': 1292, '...': 1293, '...': 1294, '...': 1295, '...': 1296, '...': 1297, '...': 1298, '...': 1299, '...': 1300, '...': 1301, '...': 1302, '...': 1303, '...': 1304, '...': 1305, '...': 1306, '...': 1307, '...': 1308, '...': 1309, '...': 1310, '...': 1311, '...': 1312, '...': 1313, '...': 1314, '...': 1315, '...': 1316, '...': 1317, '...': 1318, '...': 1319, '...': 1320, '...': 1321, '...': 1322, '...': 1323, '...': 1324, '...': 1325, '...': 1326, '...': 1327, '...': 1328, '...': 1329, '...': 1330, '...': 1331, '...': 1332, '...': 1333, '...': 1334, '...': 1335, '...': 1336, '...': 1337, '...': 1338, '...': 1339, '...': 1340, '...': 1341, '...': 1342, '...': 1343, '...': 1344, '...': 1345, '...': 1346, '...': 1347, '...': 1348, '...': 1349, '...': 1350, '...': 1351, '...': 1352, '...': 1353, '...': 1354, '...': 1355, '...': 1356, '...': 1357, '...': 1358, '...': 1359, '...': 1360, '...': 1361, '...': 1362, '...': 1363, '...': 1364, '...': 1365, '...': 1366, '...': 1367, '...': 1368, '...': 1369, '...': 1370, '...': 1371, '...': 1372, '...': 1373, '...': 1374, '...': 1375, '...': 1376, '...': 1377, '...': 1378, '...': 1379, '...': 1380, '...': 1381, '...': 1382, '...': 1383, '...': 1384, '...': 1385, '...': 1386, '...': 1387, '...': 1388, '...': 1389, '...': 1390, '...': 1391, '...': 1392, '...': 1393, '...': 1394, '...': 1395, '...': 1396, '...': 1397, '...': 1398, '...': 1399, '...': 1400, '...': 1401, '...': 1402, '...': 1403, '...': 1404, '...': 1405, '...': 1406, '...': 1407, '...': 1408, '...': 1409, '...': 1410, '...': 1411, '...': 1412, '...': 1413, '...': 1414, '...': 1415, '...': 1416, '...': 1417, '...': 1418, '...': 1419, '...': 1420, '...': 1421, '...': 1422, '...': 1423, '...': 1424, '...': 1425, '...': 1426, '...': 1427, '...': 1428, '...': 1429, '...': 1430, '...': 1431, '...': 1432, '...': 1433, '...': 1434, '...': 1435, '...': 1436, '...': 1437, '...': 1438, '...':
```



```
In [14]: #CREATE BAG OF WORDS FROM THE TOKENS FILE

tokens = [simple_preprocess(sentence, deacc = True) for sentence in open('F:\RottenTomato.csv') ]

gensim_dictionary = corpora.Dictionary ()
gensim_corpus = [gensim_dictionary.doc2bow(token, allow_update=True) for token in tokens]

word_frequencies = [[(gensim_dictionary [id], frequency)for id, frequency in couple ] for couple in gensim_corpus]

print (word_frequencies)

IOPub data rate exceeded.
The notebook server will temporarily stop sending output
to the client in order to avoid crashing it.
To change this limit, set the config variable
`--NotebookApp.iopub_data_rate_limit`.

Current values:
NotebookApp.iopub_data_rate_limit=10000000.0 (bytes/sec)
NotebookApp.rate_limit_window=3.0 (secs)
```

Fig.08

```
In [30]: # Create embedding model from text

review_model = gensim.models.Word2Vec(reviews, min_count=1)
```

Fig.09

```
In [31]: #Explore NEGATIVE patterns in text model - (most similar words)
#Source of adjectives is (https://descriptivewords.org/wp-content/uploads/descriptivewords-for-movies.pdf)

review_model.wv.most_similar(['wacky','stupid','silly','average','boring','bland', 'disappointing',
                             'distasteful', 'tiresome', 'moronic', 'predictable', 'predictable', 'silly', 'stupid', 'uninterest:

Out[31]: [('ridiculous', 0.8540648818016052),
 ('implausible', 0.8319385051727295),
 ('tasteless', 0.8295560479164124),
 ('preposterous', 0.8243200778961182),
 ('unoriginal', 0.8195405602455139),
 ('dull', 0.8086024522781372),
 ('corny', 0.8064192533493042),
 ('sappy', 0.8007218837738037),
 ('inane', 0.8003897070884705),
 ('unimaginative', 0.7996317744255066)]
```

Fig.10

```

In [42]: # Create positive embedding model from text
review_model = gensim.models.Word2Vec(reviews, min_count=10)

In [43]: #Explore POSITIVE patterns in text model - ( most similar words)
#Source of adjectives is: (https://descriptivewords.org/wp-content/uploads/descriptivewords-for-movies.pdf)

review_model.wv.most_similar(['brilliant','charismatic','charming','clever','dazzling','excellent','good',
                             'exciting','imaginative','insightful','inspirational','intriguing','legendary','original'])

Out[43]: [('engaging', 0.8168711066246033),
          ('inventive', 0.7930113673210144),
          ('captivating', 0.7663605809211731),
          ('ingenious', 0.7576407194137573),
          ('innovative', 0.7524563074111938),
          ('arresting', 0.7445882558822632),
          ('enchanted', 0.7367279529571533),
          ('endearing', 0.7363201975822449),
          ('enthraling', 0.7299625277519226),
          ('entertaining', 0.7290042638778687)]

```

Fig.11

```

In [40]: #Source of adjectives is: (https://descriptivewords.org/wp-content/uploads/descriptivewords-for-movies.pdf)

adjectives = ['weak','wacky','stupid','silly','average','boring','bland','disappointing',
              'distasteful','tiresome','moronic','predictable','predictable','silly','stupid','uninteresting',
              'brilliant','charismatic','charming','clever','dazzling','excellent','good','amazing',
              'exciting','imaginative','insightful','inspirational','intriguing','legendary','original']

#Sorting alphabetically
#adjectives.sort()

#for Loop reflecting the movie genre, the adjective and the frequency

for a in adjectives:
    print("godfather", a, review_model.wv.similarity('godfather', a))

for a in adjectives:
    print('goodfellas', a, review_model.wv.similarity('goodfellas', a))

for a in adjectives:
    print('scarface', a, review_model.wv.similarity('scarface', a))

godfather weak -0.09054291
godfather wacky 0.058230504
godfather stupid 0.018192202
godfather silly -0.017721623
godfather average 0.13879102
godfather boring 0.073217824
godfather bland 0.07282665
godfather disappointing 0.16790351
godfather distasteful -0.028919348
godfather tiresome -0.06346463
godfather moronic -0.0016719922
godfather predictable -0.012748338
godfather predictable -0.012748338
godfather silly -0.017721623
godfather stupid 0.018192202
godfather uninteresting -0.039034475

```

Fig.12

```

In [ ]: #Save the new model

review_model2 = gensim.models.Word2Vec(reviews, min_count=1)

review_model2.wv.save_word2vec_format('F:/RottenTomato02.csv')

```

Fig.13

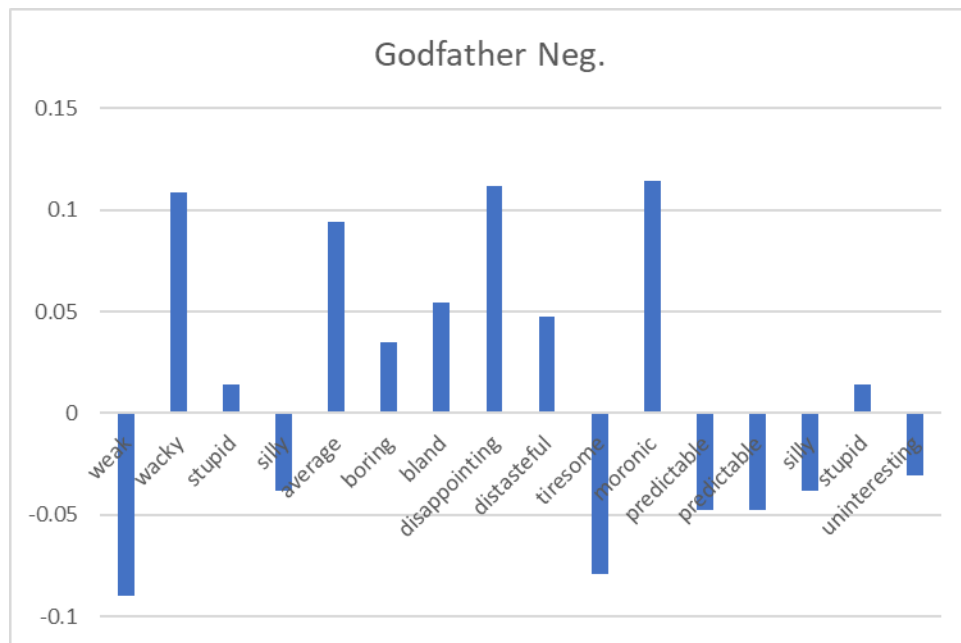


Fig.01 (Godfather movie negative adjective similarity)

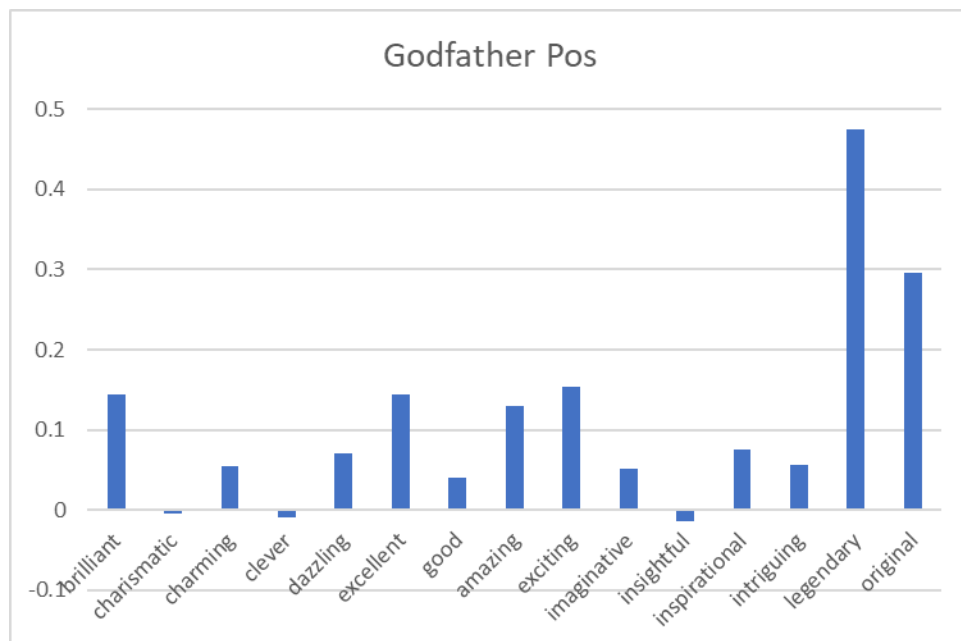


Fig.02 (Godfather movie positive adjective similarity)

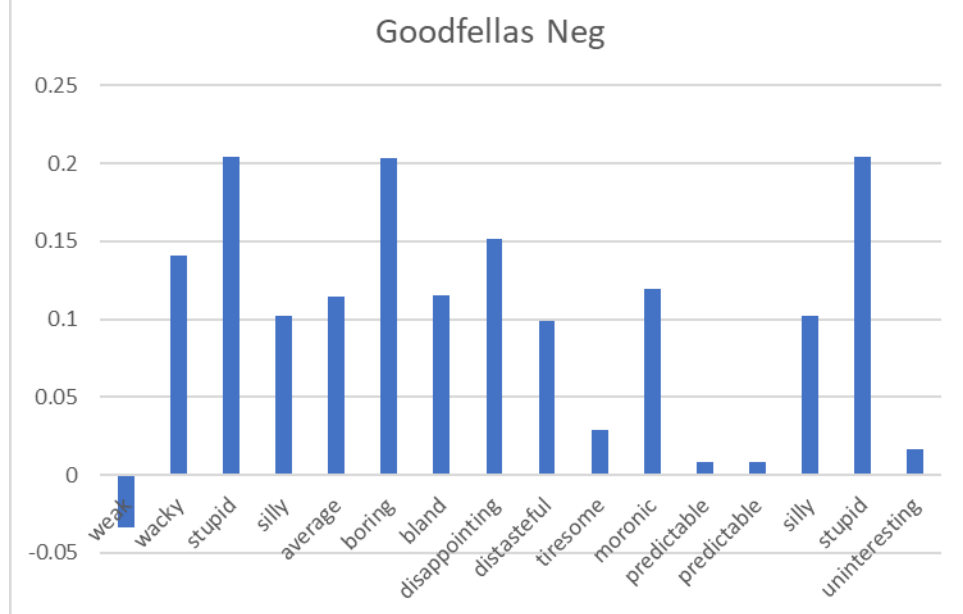


Fig.03 (Goodfellas movie negative adjective similarity)

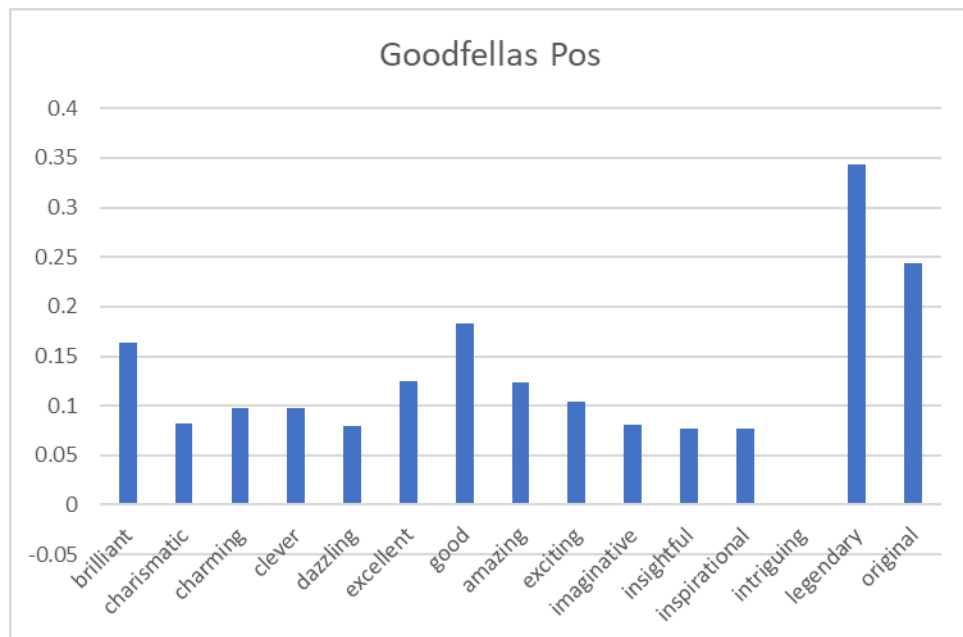


Fig.04 (Goodfellas movie positive adjective similarity)

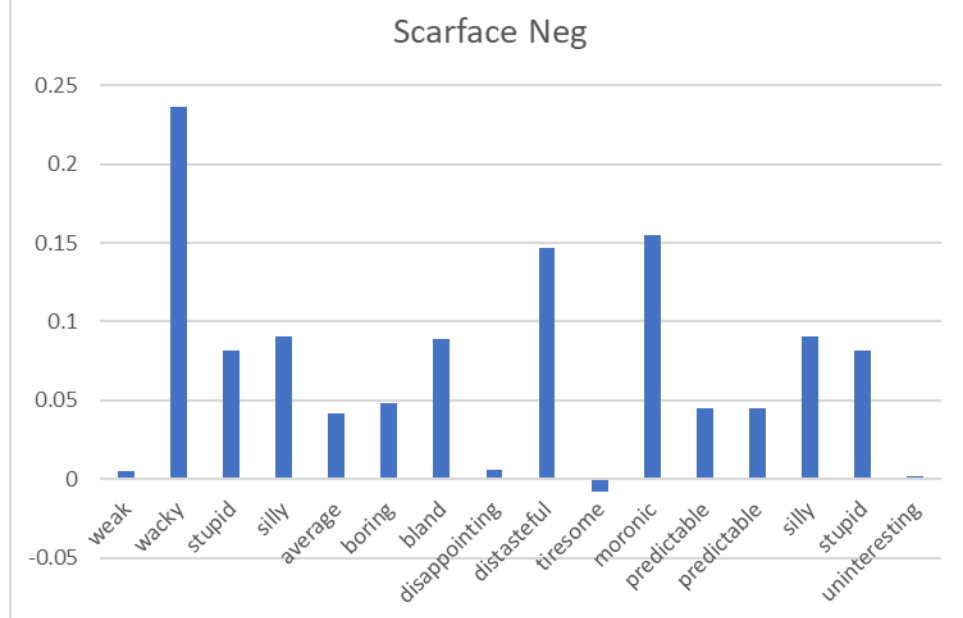


Fig.05 (Scarface movie Negative adjective similarity)

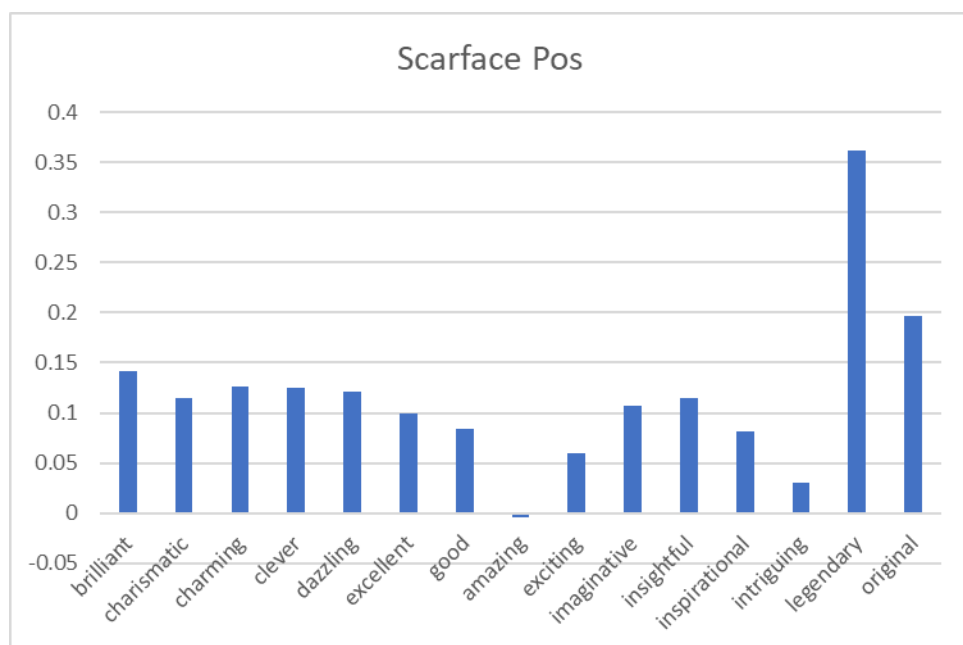


Fig.06 (Scarface movie Negative adjective similarity)