

Σύστημα συστάσεων ταινιών με χρήση αλγορίθμου
σημασιολογικής εξόρυξης

Submitted by:

Δημήτριος Τοζακίδης, **02744**



Τμήμα Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών

Πανεπιστήμιο Θεσσαλίας Βόλος, Ελλάδα

2022 - 2023

Περίληψη

Αυτό το αρχείο αποτελεί μια αναφορά ενός συστήματος συστάσεων ταινιών, με βάση τον αλγόριθμο Latent Dirichlet allocation (LDA).

Εισαγωγή

Τα συστήματα συστάσεων χρησιμεύουν παντού στο διαδίκτυο, σε διάφορες πτυχές του, όπως τα κοινωνικά δίκτυα, ιστότοποι βιβλιοθηκών, ιστότοποι προβολής ταινιών και στα συστήματα διαφημίσεων. Υπάρχουν διάφορες προσεγγίσεις για την εφαρμογή συστημάτων συστάσεων. Η LDA είναι μια δημοφιλής τεχνική στο **topic modeling**. Κυρίως χρησιμοποιείται για συστήματα συστάσεων βιβλιογραφικών εργασιών και βιβλίων. Ο λόγος είναι ότι ο αλγόριθμος αυτός βασίζεται στο πλήθος και την επανάληψη των λέξεων ενός κειμένου, ώστε να προσφέρει ικανοποιητικά αποτελέσματα. Στην αναφορά αυτή σκοπός μας είναι να δείξουμε πώς αυτή η τεχνική-αλγόριθμος μπορεί να εφαρμοστεί στον συγκεκριμένο τομέα και να μελετήσουμε πώς αποδίδει.

Εισαγωγή και καθαρισμός των εγγράφων

Ξεκινάμε βρίσκοντας τα κατάλληλα δεδομένα για το **project** μας. Το **dataframe** που επιλέξαμε αποτελείται από 14828 γραμμές και 6 στήλες. Κάνουμε **drop** την 5η στήλη διότι δεν μας χρησιμεύει. Στην συνέχεια επιλέγουμε ένα **sample 200** ταινιών ώστε να χτίσουμε το μοντέλο μας βάση αυτών. Ο καθαρισμός δεδομένων είναι απολύτως κρίσιμος για τη δημιουργία ενός χρήσιμου μοντέλου. Τα δεδομένα που μας ενδιαφέρουν για την δημιουργία του μοντέλου είναι η σύνοψη της κάθε ταινίας που χρησιμοποιούμε. Τα τρία παρακάτω βήματα είναι κοινά στις περισσότερες μεθόδους επεξεργασίας φυσικής γλώσσας:

- **Tokenizing**
- **Stopping**
- **Stemming**

Το **tokenization** τμηματοποιεί ένα έγγραφο στα ατομικά του στοιχεία. Σε αυτή την περίπτωση, μας ενδιαφέρει κάθε λέξη να αριθμηθεί. Το **tokenization** μπορεί να πραγματοποιηθεί με πολλούς τρόπους. Εμείς με την χρήση του **nlk.tokenize.simple** αποφεύ-

γουμε προβλήματα όπως το “don’t” να διαβαστεί ως δύο **tokens** “don” και “t”. Ορισμένα μέρη της αγγλικής ομιλίας, όπως οι σύνδεσμοι “for”, “or” ή η λέξη “the” δεν έχουν νόημα για ένα θέμα του μοντέλου μας. Αυτοί οι όροι ονομάζονται **stop words** και πρέπει να αφαιρεθούν από τη λίστα των **tokens** μας. Στην περίπτωση μας, χρησιμοποιούμε το πακέτο **stopwords** από την **Pypi**, μια σχετικά συντηρητική λίστα. Μπορούμε να καλέσουμε τη **getstopwords()** για να δημιουργήσουμε μια λίστα με **stop words**. Οι **Stemming words** βασίζονται στην τεχνική **NLP** για τη μείωση τοπικά όμοιων λέξεων στη ρίζα τους. Για παράδειγμα, “**stemming**,” “**stemmer**,” “**stemmed**,” όλα έχουν παρόμοιες σημασίες. Το **stemming** μειώνει αυτούς τους όρους σε “**stem**”. Αυτό είναι σημαντικό για τη μοντελοποίηση θεμάτων, η οποία διαφορετικά θα έβλεπε αυτούς τους όρους ως ξεχωριστές οντότητες και θα μείωνε τη σημασία τους στο μοντέλο. Ο αλγόριθμος **Porter stemming** είναι η πιο ευρέως χρησιμοποιούμενη μέθοδος για αυτό και την επιλέγουμε.

In [120]:

	1	movies_clean				
		imdb_id	title	plot_synopsis	tags	synopsis_source
	1305	tt0964587	St. Trinian's	Annabelle Fritton (Talulah Riley), an uptight ...	comedy	wikipedia
	3263	tt0238546	Queen of the Damned	The vampire Lestat is awakened from decades of...	paranormal, revenge, gothic, murder, flashback	wikipedia
	12892	tt0378918	True Crime: Streets of LA	Note: This plot synopsis details the "good end...	revenge, murder, violence	wikipedia
	10498	tt1259520	Bon appétit	Daniel (Unax Ugalde) is a young and ambitious ...	romantic	wikipedia
	9008	tt0297162	Half Past Dead	In San Francisco, Sasha Petrosevitch (Steven S...	violence, comedy, murder, flashback	wikipedia

	5412	tt2359024	Blue Ruin	Dwight (Macon Blair) is a beach vagrant, with ...	revenge, suspenseful, neo noir, violence	imdb
	7600	tt0332379	The School of Rock	No Vacancy, a rock band, performs at a nightcl...	cult, entertaining	wikipedia
	14530	tt0058921	Arzoo	Gopal (Rajendra Kumar) is a skiing champion. H...	romantic	wikipedia
	9913	tt0037096	Mrs. Parkington	At Christmastime in 1938, Susie Parkington, an...	revenge, flashback	wikipedia
	13853	tt0098436	The Tall Guy	The protagonist and narrator is Dexter King (G...	romantic, satire, flashback	wikipedia

200 rows x 5 columns

Figure 1: Το “καθαρισμένο” dataframe

Δημιουργία του πίνακα

Το αποτέλεσμα του σταδίου καθαρισμού μας είναι μια **tokenized**, **stopped** και **stemmed** λίστα λέξεων από ένα μόνο έγγραφο. Στην συνέχεια εφαρμόζουμε αυτή την μέθοδο σε κάθε σύνοψη από τις ταινίες που χρησιμοποιούμε. Για να δημιουργήσουμε ένα μοντέλο **LDA**, πρέπει να κατανοήσουμε πόσο συχνά εμφανίζεται κάθε όρος σε κάθε έγγραφο. Για να γίνει αυτό, πρέπει να κατασκευάσουμε ένα έγγραφο με ένα πακέτο που ονομάζεται **gensim**. Η συνάρτηση **Dictionary()** διασχίζει το έγγραφο εκχωρώντας ένα μοναδικό ακέραιο σε κάθε μοναδικό **token**, ενώ παράλληλα συλλέγει πλήθος λέξεων και σχετικά

στατιστικά στοιχεία.

```
In [64]: 1 print(dictionary.token2id)

{'19th': 0, 'abandon': 1, 'abbi': 2, 'absenc': 3, 'acquisit': 4, 'actual': 5, 'agre': 6, 'aid': 7, 'air': 8, 'alfonsi': 9, 'ali': 10, 'alreadi': 11, 'also': 12, 'anderson': 13, 'anticip': 14, 'apart': 15, 'appar': 16, 'appear': 17, 'arm': 18, 'around': 19, 'arriv': 20, 'ask': 21, 'assail': 22, 'attack': 23, 'attempt': 24, 'attract': 25, 'await': 26, 'awaken': 27, 'away': 28, 'back': 29, 'basement': 30, 'battl': 31, 'bed': 32, 'beg': 33, 'begin': 34, 'behead': 35, 'behind': 36, 'believ': 37, 'belong': 38, 'beset': 39, 'best': 40, 'bit': 41, 'bite': 42, 'blade': 43, 'blood': 44, 'bode': 45, 'bodi': 46, 'bori': 47, 'branch': 48, 'break': 49, 'breakup': 50, 'brother': 51, 'bruise': 52, 'buri': 53, 'burial': 54, 'busi': 55, 'butcher': 56, 'buzz': 57, 'cadav': 58, 'call': 59, 'caller': 60, 'calm': 61, 'camera': 62, 'can': 63, 'care': 64, 'case': 65, 'cathedr': 66, 'caus': 67, 'centuri': 68, 'certain': 69, 'chang': 70, 'chester': 71, 'child': 72, 'close': 73, 'come': 74, 'commit': 75, 'conciern': 76, 'conclud': 77, 'confess': 78, 'confirm': 79, 'confus': 80, 'continu': 81, 'corp': 82, 'cottage': 83, 'count': 84, 'cours': 85, 'creep': 86, 'crew': 87, 'crewmen': 88, 'curs': 89, 'd': 90, 'dagger': 91, 'daughter': 92, 'dawn': 93, 'day': 94, 'de': 95, 'dead': 96, 'death': 97, 'decid': 98, 'decis': 99, 'demeanor': 100, 'dialina': 101, 'diamond': 102, 'die': 103, 'difficult': 104, 'discov': 105, 'distinct': 106, 'distract': 107, 'distress': 108, 'doctor': 109, 'door': 110, 'doubt': 111, 'drain': 112, 'dread': 113, 'dress': 114, 'drip': 115, 'drop': 116, 'duti': 117, 'eagerli': 118, 'earlier': 119, 'east': 120, 'elabor': 121, 'elderli': 122, 'embrac': 123, 'end': 124, 'england': 125, 'entranc': 126, 'escap': 127, 'estrang': 128, 'even': 129, 'event': 130, 'ex': 131, 'examin': 132, 'explain': 133, 'face': 134, 'fail': 135, 'fake': 136, 'famili': 137, 'father': 138, 'fear': 139, 'feed': 140, 'felt': 141, 'final': 142, 'find': 143, 'finger': 144, 'finish': 145, 'fit': 146, 'five': 147, 'flat': 148, 'flee': 149, 'fli': 150, 'floor': 151, 'follow': 152, 'forc': 153, 'forest': 154, 'forgiv': 155, 'former': 156, 'frank': 157, 'friend': 158, 'fright': 159, 'front': 160, 'gasp': 161, 'get': 162, 'ghost': 163, 'giorgio': 164, 'girl': 165, 'give': 166, 'glass': 167, 'glauco': 168, 'go': 169, 'gold': 170, 'good': 171, 'govern': 172, 'grace': 173, 'grand': 174, 'grasp': 175, 'great': 176, 'green': 177, 'grip': 178, 'grove': 179, 'grow': 180, 'guard': 181, 'guess': 182, 'guest': 183, 'guide': 184, 'gun': 185, 'h': 186, 'hail': 187, 'hand': 188, 'hang': 189, 'harm': 190, 'has': 191, 'have': 192, 'head': 193, 'hear': 194, 'heart': 195, 'heat': 196, 'help': 197, 'here': 198, 'high': 199, 'him': 200, 'his': 201, 'hold': 202, 'hole': 203, 'home': 204, 'honest': 205, 'horse': 206, 'hospital': 207, 'host': 208, 'hot': 209, 'house': 210, 'how': 211, 'hug': 212, 'human': 213, 'hundred': 214, 'hurt': 215, 'husband': 216, 'hundred': 217, 'hundred': 218, 'hundred': 219, 'hundred': 220, 'hundred': 221, 'hundred': 222, 'hundred': 223, 'hundred': 224, 'hundred': 225, 'hundred': 226, 'hundred': 227, 'hundred': 228, 'hundred': 229, 'hundred': 230, 'hundred': 231, 'hundred': 232, 'hundred': 233, 'hundred': 234, 'hundred': 235, 'hundred': 236, 'hundred': 237, 'hundred': 238, 'hundred': 239, 'hundred': 240, 'hundred': 241, 'hundred': 242, 'hundred': 243, 'hundred': 244, 'hundred': 245, 'hundred': 246, 'hundred': 247, 'hundred': 248, 'hundred': 249, 'hundred': 250, 'hundred': 251, 'hundred': 252, 'hundred': 253, 'hundred': 254, 'hundred': 255, 'hundred': 256, 'hundred': 257, 'hundred': 258, 'hundred': 259, 'hundred': 260, 'hundred': 261, 'hundred': 262, 'hundred': 263, 'hundred': 264, 'hundred': 265, 'hundred': 266, 'hundred': 267, 'hundred': 268, 'hundred': 269, 'hundred': 270, 'hundred': 271, 'hundred': 272, 'hundred': 273, 'hundred': 274, 'hundred': 275, 'hundred': 276, 'hundred': 277, 'hundred': 278, 'hundred': 279, 'hundred': 280, 'hundred': 281, 'hundred': 282, 'hundred': 283, 'hundred': 284, 'hundred': 285, 'hundred': 286, 'hundred': 287, 'hundred': 288, 'hundred': 289, 'hundred': 290, 'hundred': 291, 'hundred': 292, 'hundred': 293, 'hundred': 294, 'hundred': 295, 'hundred': 296, 'hundred': 297, 'hundred': 298, 'hundred': 299, 'hundred': 300, 'hundred': 301, 'hundred': 302, 'hundred': 303, 'hundred': 304, 'hundred': 305, 'hundred': 306, 'hundred': 307, 'hundred': 308, 'hundred': 309, 'hundred': 310, 'hundred': 311, 'hundred': 312, 'hundred': 313, 'hundred': 314, 'hundred': 315, 'hundred': 316, 'hundred': 317, 'hundred': 318, 'hundred': 319, 'hundred': 320, 'hundred': 321, 'hundred': 322, 'hundred': 323, 'hundred': 324, 'hundred': 325, 'hundred': 326, 'hundred': 327, 'hundred': 328, 'hundred': 329, 'hundred': 330, 'hundred': 331, 'hundred': 332, 'hundred': 333, 'hundred': 334, 'hundred': 335, 'hundred': 336, 'hundred': 337, 'hundred': 338, 'hundred': 339, 'hundred': 340, 'hundred': 341, 'hundred': 342, 'hundred': 343, 'hundred': 344, 'hundred': 345, 'hundred': 346, 'hundred': 347, 'hundred': 348, 'hundred': 349, 'hundred': 350, 'hundred': 351, 'hundred': 352, 'hundred': 353, 'hundred': 354, 'hundred': 355, 'hundred': 356, 'hundred': 357, 'hundred': 358, 'hundred': 359, 'hundred': 360, 'hundred': 361, 'hundred': 362, 'hundred': 363, 'hundred': 364, 'hundred': 365, 'hundred': 366, 'hundred': 367, 'hundred': 368, 'hundred': 369, 'hundred': 370, 'hundred': 371, 'hundred': 372, 'hundred': 373, 'hundred': 374, 'hundred': 375, 'hundred': 376, 'hundred': 377, 'hundred': 378, 'hundred': 379, 'hundred': 380, 'hundred': 381, 'hundred': 382, 'hundred': 383, 'hundred': 384, 'hundred': 385, 'hundred': 386, 'hundred': 387, 'hundred': 388, 'hundred': 389, 'hundred': 390, 'hundred': 391, 'hundred': 392, 'hundred': 393, 'hundred': 394, 'hundred': 395, 'hundred': 396, 'hundred': 397, 'hundred': 398, 'hundred': 399, 'hundred': 400, 'hundred': 401, 'hundred': 402, 'hundred': 403, 'hundred': 404, 'hundred': 405, 'hundred': 406, 'hundred': 407, 'hundred': 408, 'hundred': 409, 'hundred': 410, 'hundred': 411, 'hundred': 412, 'hundred': 413, 'hundred': 414, 'hundred': 415, 'hundred': 416, 'hundred': 417, 'hundred': 418, 'hundred': 419, 'hundred': 420, 'hundred': 421, 'hundred': 422, 'hundred': 423, 'hundred': 424, 'hundred': 425, 'hundred': 426, 'hundred': 427, 'hundred': 428, 'hundred': 429, 'hundred': 430, 'hundred': 431, 'hundred': 432, 'hundred': 433, 'hundred': 434, 'hundred': 435, 'hundred': 436, 'hundred': 437, 'hundred': 438, 'hundred': 439, 'hundred': 440, 'hundred': 441, 'hundred': 442, 'hundred': 443, 'hundred': 444, 'hundred': 445, 'hundred': 446, 'hundred': 447, 'hundred': 448, 'hundred': 449, 'hundred': 450, 'hundred': 451, 'hundred': 452, 'hundred': 453, 'hundred': 454, 'hundred': 455, 'hundred': 456, 'hundred': 457, 'hundred': 458, 'hundred': 459, 'hundred': 460, 'hundred': 461, 'hundred': 462, 'hundred': 463, 'hundred': 464, 'hundred': 465, 'hundred': 466, 'hundred': 467, 'hundred': 468, 'hundred': 469, 'hundred': 470, 'hundred': 471, 'hundred': 472, 'hundred': 473, 'hundred': 474, 'hundred': 475, 'hundred': 476, 'hundred': 477, 'hundred': 478, 'hundred': 479, 'hundred': 480, 'hundred': 481, 'hundred': 482, 'hundred': 483, 'hundred': 484, 'hundred': 485, 'hundred': 486, 'hundred': 487, 'hundred': 488, 'hundred': 489, 'hundred': 490, 'hundred': 491, 'hundred': 492, 'hundred': 493, 'hundred': 494, 'hundred': 495, 'hundred': 496, 'hundred': 497, 'hundred': 498, 'hundred': 499, 'hundred': 500, 'hundred': 501, 'hundred': 502, 'hundred': 503, 'hundred': 504, 'hundred': 505, 'hundred': 506, 'hundred': 507, 'hundred': 508, 'hundred': 509, 'hundred': 510, 'hundred': 511, 'hundred': 512, 'hundred': 513, 'hundred': 514, 'hundred': 515, 'hundred': 516, 'hundred': 517, 'hundred': 518, 'hundred': 519, 'hundred': 520, 'hundred': 521, 'hundred': 522, 'hundred': 523, 'hundred': 524, 'hundred': 525, 'hundred': 526, 'hundred': 527, 'hundred': 528, 'hundred': 529, 'hundred': 530, 'hundred': 531, 'hundred': 532, 'hundred': 533, 'hundred': 534, 'hundred': 535, 'hundred': 536, 'hundred': 537, 'hundred': 538, 'hundred': 539, 'hundred': 540, 'hundred': 541, 'hundred': 542, 'hundred': 543, 'hundred': 544, 'hundred': 545, 'hundred': 546, 'hundred': 547, 'hundred': 548, 'hundred': 549, 'hundred': 550, 'hundred': 551, 'hundred': 552, 'hundred': 553, 'hundred': 554, 'hundred': 555, 'hundred': 556, 'hundred': 557, 'hundred': 558, 'hundred': 559, 'hundred': 560, 'hundred': 561, 'hundred': 562, 'hundred': 563, 'hundred': 564, 'hundred': 565, 'hundred': 566, 'hundred': 567, 'hundred': 568, 'hundred': 569, 'hundred': 570, 'hundred': 571, 'hundred': 572, 'hundred': 573, 'hundred': 574, 'hundred': 575, 'hundred': 576, 'hundred': 577, 'hundred': 578, 'hundred': 579, 'hundred': 580, 'hundred': 581, 'hundred': 582, 'hundred': 583, 'hundred': 584, 'hundred': 585, 'hundred': 586, 'hundred': 587, 'hundred': 588, 'hundred': 589, 'hundred': 590, 'hundred': 591, 'hundred': 592, 'hundred': 593, 'hundred': 594, 'hundred': 595, 'hundred': 596, 'hundred': 597, 'hundred': 598, 'hundred': 599, 'hundred': 600, 'hundred': 601, 'hundred': 602, 'hundred': 603, 'hundred': 604, 'hundred': 605, 'hundred': 606, 'hundred': 607, 'hundred': 608, 'hundred': 609, 'hundred': 610, 'hundred': 611, 'hundred': 612, 'hundred': 613, 'hundred': 614, 'hundred': 615, 'hundred': 616, 'hundred': 617, 'hundred': 618, 'hundred': 619, 'hundred': 620, 'hundred': 621, 'hundred': 622, 'hundred': 623, 'hundred': 624, 'hundred': 625, 'hundred': 626, 'hundred': 627, 'hundred': 628, 'hundred': 629, 'hundred': 630, 'hundred': 631, 'hundred': 632, 'hundred': 633, 'hundred': 634, 'hundred': 635, 'hundred': 636, 'hundred': 637, 'hundred': 638, 'hundred': 639, 'hundred': 640, 'hundred': 641, 'hundred': 642, 'hundred': 643, 'hundred': 644, 'hundred': 645, 'hundred': 646, 'hundred': 647, 'hundred': 648, 'hundred': 649, 'hundred': 650, 'hundred': 651, 'hundred': 652, 'hundred': 653, 'hundred': 654, 'hundred': 655, 'hundred': 656, 'hundred': 657, 'hundred': 658, 'hundred': 659, 'hundred': 660, 'hundred': 661, 'hundred': 662, 'hundred': 663, 'hundred': 664, 'hundred': 665, 'hundred': 666, 'hundred': 667, 'hundred': 668, 'hundred': 669, 'hundred': 670, 'hundred': 671, 'hundred': 672, 'hundred': 673, 'hundred': 674, 'hundred': 675, 'hundred': 676, 'hundred': 677, 'hundred': 678, 'hundred': 679, 'hundred': 680, 'hundred': 681, 'hundred': 682, 'hundred': 683, 'hundred': 684, 'hundred': 685, 'hundred': 686, 'hundred': 687, 'hundred': 688, 'hundred': 689, 'hundred': 690, 'hundred': 691, 'hundred': 692, 'hundred': 693, 'hundred': 694, 'hundred': 695, 'hundred': 696, 'hundred': 697, 'hundred': 698, 'hundred': 699, 'hundred': 700, 'hundred': 701, 'hundred': 702, 'hundred': 703, 'hundred': 704, 'hundred': 705, 'hundred': 706, 'hundred': 707, 'hundred': 708, 'hundred': 709, 'hundred': 710, 'hundred': 711, 'hundred': 712, 'hundred': 713, 'hundred': 714, 'hundred': 715, 'hundred': 716, 'hundred': 717, 'hundred': 718, 'hundred': 719, 'hundred': 720, 'hundred': 721, 'hundred': 722, 'hundred': 723, 'hundred': 724, 'hundred': 725, 'hundred': 726, 'hundred': 727, 'hundred': 728, 'hundred': 729, 'hundred': 730, 'hundred': 731, 'hundred': 732, 'hundred': 733, 'hundred': 734, 'hundred': 735, 'hundred': 736, 'hundred': 737, 'hundred': 738, 'hundred': 739, 'hundred': 740, 'hundred': 741, 'hundred': 742, 'hundred': 743, 'hundred': 744, 'hundred': 745, 'hundred': 746, 'hundred': 747, 'hundred': 748, 'hundred': 749, 'hundred': 750, 'hundred': 751, 'hundred': 752, 'hundred': 753, 'hundred': 754, 'hundred': 755, 'hundred': 756, 'hundred': 757, 'hundred': 758, 'hundred': 759, 'hundred': 760, 'hundred': 761, 'hundred': 762, 'hundred': 763, 'hundred': 764, 'hundred': 765, 'hundred': 766, 'hundred': 767, 'hundred': 768, 'hundred': 769, 'hundred': 770, 'hundred': 771, 'hundred': 772, 'hundred': 773, 'hundred': 774, 'hundred': 775, 'hundred': 776, 'hundred': 777, 'hundred': 778, 'hundred': 779, 'hundred': 780, 'hundred': 781, 'hundred': 782, 'hundred': 783, 'hundred': 784, 'hundred': 785, 'hundred': 786, 'hundred': 787, 'hundred': 788, 'hundred': 789, 'hundred': 790, 'hundred': 791, 'hundred': 792, 'hundred': 793, 'hundred': 794, 'hundred': 795, 'hundred': 796, 'hundred': 797, 'hundred': 798, 'hundred': 799, 'hundred': 800, 'hundred': 801, 'hundred': 802, 'hundred': 803, 'hundred': 804, 'hundred': 805, 'hundred': 806, 'hundred': 807, 'hundred': 808, 'hundred': 809, 'hundred': 810, 'hundred': 811, 'hundred': 812, 'hundred': 813, 'hundred': 814, 'hundred': 815, 'hundred': 816, 'hundred': 817, 'hundred': 818, 'hundred': 819, 'hundred': 820, 'hundred': 821, 'hundred': 822, 'hundred': 823, 'hundred': 824, 'hundred': 825, 'hundred': 826, 'hundred': 827, 'hundred': 828, 'hundred': 829, 'hundred': 830, 'hundred': 831, 'hundred': 832, 'hundred': 833, 'hundred': 834, 'hundred': 835, 'hundred': 836, 'hundred': 837, 'hundred': 838, 'hundred': 839, 'hundred': 840, 'hundred': 841, 'hundred': 842, 'hundred': 843, 'hundred': 844, 'hundred': 845, 'hundred': 846, 'hundred': 847, 'hundred': 848, 'hundred': 849, 'hundred': 850, 'hundred': 851, 'hundred': 852, 'hundred': 853, 'hundred': 854, 'hundred': 855, 'hundred': 856, 'hundred': 857, 'hundred': 858, 'hundred': 859, 'hundred': 860, 'hundred': 861, 'hundred': 862, 'hundred': 863, 'hundred': 864, 'hundred': 865, 'hundred': 866, 'hundred': 867, 'hundred': 868, 'hundred': 869, 'hundred': 870, 'hundred': 871, 'hundred': 872, 'hundred': 873, 'hundred': 874, 'hundred': 875, 'hundred': 876, 'hundred': 877, 'hundred': 878, 'hundred': 879, 'hundred': 880, 'hundred': 881, 'hundred': 882, 'hundred': 883, 'hundred': 884, 'hundred': 885, 'hundred': 886, 'hundred': 887, 'hundred': 888, 'hundred': 889, 'hundred': 890, 'hundred': 891, 'hundred': 892, 'hundred': 893, 'hundred': 894, 'hundred': 895, 'hundred': 896, 'hundred': 897, 'hundred': 898, 'hundred': 899, 'hundred': 900, 'hundred': 901, 'hundred': 902, 'hundred': 903, 'hundred': 904, 'hundred': 905, 'hundred': 906, 'hundred': 907, 'hundred': 908, 'hundred': 909, 'hundred': 910, 'hundred': 911, 'hundred': 912, 'hundred': 913, 'hundred': 914, 'hundred': 915, 'hundred': 916, 'hundred': 917, 'hundred': 918, 'hundred': 919, 'hundred': 920, 'hundred': 921, 'hundred': 922, 'hundred': 923, 'hundred': 924, 'hundred': 925, 'hundred': 926, 'hundred': 927, 'hundred': 928, 'hundred': 929, 'hundred': 930, 'hundred': 931, 'hundred': 932, 'hundred': 933, 'hundred': 934, 'hundred': 935, 'hundred': 936, 'hundred': 937, 'hundred': 938, 'hundred': 939, 'hundred': 940, 'hundred': 941, 'hundred': 942, 'hundred': 943, 'hundred': 944, 'hundred': 945, 'hundred': 946, 'hundred': 947, 'hundred': 948, 'hundred': 949, 'hundred': 950, 'hundred': 951, 'hundred': 952, 'hundred': 953, 'hundred': 954, 'hundred': 955, 'hundred': 956, 'hundred': 957, 'hundred': 958, 'hundred': 959, 'hundred': 960, 'hundred': 961, 'hundred': 962, 'hundred': 963, 'hundred': 964, 'hundred': 965, 'hundred': 966, 'hundred': 967, 'hundred': 968, 'hundred': 969, 'hundred': 970, 'hundred': 971, 'hundred': 972, 'hundred': 973, 'hundred': 974, 'hundred': 975, 'hundred': 976, 'hundred': 977, 'hundred': 978, 'hundred': 979, 'hundred': 980, 'hundred': 981, 'hundred': 982, 'hundred': 983, 'hundred': 984, 'hundred': 985, 'hundred': 986, 'hundred': 987, 'hundred': 988, 'hundred': 989, 'hundred': 990, 'hundred': 991, 'hundred': 992, 'hundred': 993, 'hundred': 994, 'hundred': 995, 'hundred': 996, 'hundred': 997, 'hundred': 998, 'hundred': 999, 'hundred': 1000}
```

Figure 2: Το "λεξικό" με τις λέξεις μας

Η συνάρτηση `doc2bow()` μετατρέπει το `Dictionary()` σε τσάντα με λέξεις. Το αποτέλεσμα, `corpus`, είναι μια λίστα `vectors` ίση με τον αριθμό των εγγράφων. Σε κάθε `vector` εγγράφου υπάρχει μια σειρά από πλειάδες. Για παράδειγμα, το `print(corpus[0])` έχει ως αποτέλεσμα τα εξής:

```
In [66]: 1 print(corpus[0])

[(0, 1), (1, 2), (2, 1), (3, 1), (4, 1), (5, 1), (6, 2), (7, 2), (8, 1), (9, 1), (10, 2), (11, 1), (12, 2), (13, 1), (14, 1), (15, 4), (16, 1), (17, 3), (18, 1), (19, 1), (20, 3), (21, 1), (22, 2), (23, 1), (24, 1), (25, 2), (26, 1), (27, 3), (28, 1), (29, 2), (30, 1), (31, 1), (32, 5), (33, 4), (34, 2), (35, 2), (36, 1), (37, 1), (38, 1), (39, 1), (40, 1), (41, 1), (42, 2), (43, 1), (44, 3), (45, 1), (46, 4), (47, 3), (48, 1), (49, 2), (50, 1), (51, 2), (52, 1), (53, 1), (54, 1), (55, 1), (56, 1), (57, 1), (58, 1), (59, 8), (60, 2), (61, 1), (62, 1), (63, 1), (64, 2), (65, 1), (66, 1), (67, 1), (68, 1), (69, 3), (70, 1), (71, 6), (72, 2), (73, 1), (74, 3), (75, 1), (76, 2), (77, 2), (78, 1), (79, 1), (80, 1), (81, 2), (82, 5), (83, 5), (84, 2), (85, 1), (86, 1), (87, 1), (88, 1), (89, 1), (90, 1), (91, 1), (92, 1), (93, 1), (94, 1), (95, 1), (96, 1), (97, 2), (98, 1), (99, 1), (100, 1), (101, 1), (102, 1), (103, 1), (104, 1), (105, 3), (106, 1), (107, 1), (108, 1), (109, 1), (110, 1), (111, 1), (112, 1), (113, 1), (114, 1), (115, 1), (116, 3), (117, 1), (118, 1), (119, 1), (120, 1), (121, 1), (122, 1), (123, 1), (124, 1), (125, 1), (126, 1), (127, 2), (128, 1), (129, 2), (130, 1), (131, 1), (132, 1), (133, 3), (134, 2), (135, 1), (136, 1), (137, 7), (138, 3), (139, 4), (140, 2), (141, 1), (142, 1), (143, 2), (144, 3), (145, 1), (146, 1), (147, 1), (148, 1), (149, 2), (150, 3), (151, 1), (152, 1), (153, 1), (154, 1), (155, 1), (156, 1), (157, 6), (158, 1), (159, 2), (160, 1), (161, 1), (162, 2), (163, 1), (164, 6), (165, 1), (166, 1), (167, 1), (168, 1), (169, 2), (170, 1), (171, 7), (172, 1), (173, 1), (174, 1), (175, 1), (176, 1), (177, 1), (178, 1), (179, 1), (180, 1), (181, 1), (182, 1), (183, 1), (184, 2), (185, 1), (186, 1), (187, 4), (188, 1), (189, 2), (190, 1), (191, 1), (192, 1), (193, 1), (194, 1), (195, 2), (196, 1), (197, 1), (198, 3), (199, 2), (200, 1), (201, 1), (202, 1), (203, 1), (204, 3), (205, 1), (206, 1), (207, 1), (208, 1), (209, 1), (210, 3), (211, 2), (212, 5), (213, 1), (214, 1), (215, 1), (216, 1), (217, 4), (218, 1), (219, 1), (220, 1), (221, 1), (222, 1), (223, 3), (224, 2), (225, 1), (226, 1), (227, 1), (228, 1), (229, 2), (230, 1), (231, 1), (232, 2), (233, 3), (234, 1), (235, 1), (236, 1), (237, 1), (238, 5), (239, 3), (240, 6), (241, 1), (242, 1), (243, 1), (244, 1), (245, 1), (246, 2), (247, 1), (248, 1), (249, 1), (250, 1), (251, 1), (252, 1), (253, 1), (254, 1), (255, 1), (256, 1), (257, 1), (258, 1), (259, 1), (260, 1), (261, 7), (262, 2), (263, 3), (264, 3), (265, 6), (266, 1), (267, 1), (268, 1), (269, 1), (270, 3), (271, 5), (272, 1), (273, 1), (274, 1), (275, 1), (276, 1), (277, 1), (278, 1), (279, 1), (280, 1), (281, 1), (282, 3), (283, 1), (284, 4), (285, 1), (286, 2), (287, 1), (288, 1), (289, 1), (290, 1), (291, 1), (292, 1), (293, 1), (294, 1), (295, 2), (296, 1), (297, 1), (298, 2), (299, 1), (300, 1), (301, 1), (302, 1), (303, 1), (304, 1), (305, 1), (306, 1), (307, 3), (308, 1), (309, 1), (310, 1), (311, 1), (312, 1), (313, 1), (314, 1), (315, 1), (316, 6), (317, 1), (318, 1), (319, 1), (320, 2), (321, 1), (322, 1), (323, 3), (324, 1), (325, 1), (326, 14), (327, 1), (328, 2), (329, 1), (330, 1), (331, 1), (332, 11), (333, 3), (334, 1), (335, 7), (336, 1), (337, 1), (338, 1), (339, 2), (340, 1), (341, 1), (342, 1), (343, 1), (344, 1), (345, 1), (346, 1), (347, 1), (348, 1), (349, 1), (350, 1), (351, 2), (352, 1), (353, 3), (354, 1), (355, 1), (356, 2), (357, 2), (358, 3), (359, 1), (360, 1), (361, 1), (362, 1), (363, 1), (364, 2), (365, 1), (366, 2), (367, 1), (368, 1), (369, 2), (370, 3), (371, 3), (372, 1), (373, 1), (374, 1), (375, 1), (376, 2), (377, 1), (378, 1), (379, 1), (380, 1), (381, 1), (382, 1), (383, 1), (384, 1), (385, 1), (386, 2), (387, 2), (388, 2), (389, 1), (390, 1), (391, 1), (392, 1), (393, 2), (394, 1), (395, 1), (396, 3), (397, 1), (398, 1), (399, 2), (400, 1), (401, 1), (402, 1), (403, 1), (404, 1), (405, 1), (406, 1), (407, 3), (408, 1), (409, 1), (410, 1), (411, 1), (412, 1), (413, 1), (414, 1), (415, 1), (416, 1), (417, 1), (418, 1), (419, 1), (420, 10), (421, 2), (422, 1),
```

Εφαρμόζοντας το **LDA** μοντέλο

Το **corpus** είναι ένας πίνακας εγγράφων και τώρα είμαστε έτοιμοι να δημιουργήσουμε ένα μοντέλο **LDA**. Η κλάση **LdaModel** ανήκει στο **gensim** και οι παράμετροι που χρησιμοποιούμε είναι οι εξής:

- **numtopics**: Απαιτείται. Ένα μοντέλο **LDA** απαιτεί από τον χρήστη να καθορίσει πόσα θέματα πρέπει να δημιουργηθούν. Μετά από δοκιμές επιλέξαμε να ζητήσουμε **30** θέματα
- **id2word**: Απαιτείται. Η κλάση **LdaModel** απαιτεί από το προηγούμενο **dictionary** μας να αντιστοιχίσει τους μοναδικούς ακεραίους σε συμβολοσειρές.
- **passes**: Προαιρετικός. Ο αριθμός των επαναλήψεων που θα κάνει το μοντέλο μέσω του **corpus**. Όσο μεγαλύτερος είναι ο αριθμός των περασμάτων, τόσο πιο ακριβές θα είναι το μοντέλο. Πολλά περάσματα μπορεί να είναι αργά σε ένα πολύ μεγάλο **corpus**.

Εδώ εκτυπώνουμε και τα **30** θέματα. Κάθε καταχώριση θα έχει τις **10** κορυφαίες λέξεις κατά σημαντικότητα, ακολουθούμενες από τη συσχέτιση του συνόλου των θεμάτων, τον μέσο όρο των βαθμολογιών κατά ζεύγη ομοιότητας λέξεων όλων των λέξεων εντός του θέματος.

```
In [121]: 1 from pprint import pprint
          2 top_topics = lda_model_final.top_topics(corpus, topn=10)

In [122]: 1 counter = 0
          2 for topic in top_topics:
          3     print('Topic {}'.format(counter))
          4     counter += 1
          5     pprint(topic)

Topic 0:
[(0.05555626, 'harrer'),
 (0.024083028, 'aufschnait'),
 (0.01983849, 'chines'),
 (0.01945063, 'lama'),
 (0.018528925, 'tibetan'),
 (0.016672892, 'dalai'),
 (0.015748698, 'tibet'),
 (0.012974823, 'lhasa'),
 (0.010197772, 'jigm'),
 (0.00754309, 'gift')],
-0.4684910551240771)
Topic 1:
[(0.02063212, 'terri'),
 (0.01857233, 'ian'),
 (0.010655121, 'howard'),
 (0.010333157, 'brother'),
 (0.008542705, 'kill'),
 (0.008378023, 'burn'),
 (0.003300105, 't')]
```

Figure 4: Τα 30 θέματα

Ας δούμε τώρα από ποια θέματα αποτελείται η ταινία **St. Trinian's** που είναι η πρώτη ταινία στο **dataframe** των **200** ταινιών:

Το υψηλότερο συσχετιζόμενο θέμα ήταν το **11**, το οποίο είναι: **'oak', 'telli', 'calvess', 'darnel', 't', 'gun', 'gump', 'say', 'ask', 'talk'**.

```
In [154]: 1 corpus_lda_model = lda_model_final[corpus]
2         for article in corpus_lda_model[0]:
3             print(article)

(11, 0.80746126)
(22, 0.19123007)
```

```
In [155]: 1 movies_clean
```

```
Out[155]:
```

	imdb_id	title	plot_synopsis	tags	synopsis_source
1305	tt0964587	St. Trinian's	Annabelle Fritton (Talulah Riley), an uptight ...	comedy	wikipedia
3263	tt0238546	Queen of the Damned	The vampire Lestat is awakened from decades of...	paranormal, revenge, gothic, murder, flashback	wikipedia
12892	tt0378918	True Crime: Streets of LA	Note: This plot synopsis details the "good end...	revenge, murder, violence	wikipedia
10498	tt1259520	Bon appétit	Daniel (Unax Ugalde) is a young and ambitious ...	romantic	wikipedia
9008	tt0297162	Half Past Dead	In San Francisco, Sasha Petroskevitch (Steven S...	violence, comedy, murder, flashback	wikipedia
...
5412	tt2359024	Blue Ruin	Dwight (Macon Blair) is a beach vagrant, with ...	revenge, suspenseful, neo noir, violence	imdb
7600	tt0332379	The School of Rock	No Vacancy, a rock band, performs at a nightcl...	cult, entertaining	wikipedia
14530	tt0058921	Arzoo	Gopal (Rajendra Kumar) is a skiing champion. H...	romantic	wikipedia
9913	tt0037096	Mrs. Parkington	At Christmastime in 1938, Susie Parkington, an...	revenge, flashback	wikipedia
13853	tt0098436	The Tall Guy	The protagonist and narrator is Dexter King (G...	romantic, satire, flashback	wikipedia

200 rows × 5 columns

Figure 5: Τα θέματα της ταινίας St. Trinian's

Ας κάνουμε έναν τελικό έλεγχο με βάση την σειρά των εγγράφων μας κατά μήκος. Θα ήταν λογικό τα μεγαλύτερα άρθρα να αποτελούνται από περισσότερα θέματα από τα μικρότερα, οπότε ας γράψουμε τη θεματική συσχέτιση της οπτικής με την υψηλότερη συσχέτιση για κάθε άρθρο και ας δούμε αν μπορούν να έχουν μια φθίνουσα σχέση.

```

In [149]: 1 from operator import itemgetter
2 index_high_corr_list = []
3
4 index_counter = 0
5
6 for article in corpus_lda_model:
7     highest_feature = max(article, key=itemgetter(1))[0]
8     max_correlation = max(article, key=itemgetter(1))[1]
9     index_high_corr_list.append([index_counter, highest_feature, max_correlation])
10    index_counter += 1
11
12 index_high_corr_list[:5]

```

```

Out[149]: [[0, 29, 0.55548346],
[1, 25, 0.99485016],
[2, 34, 0.45816225],
[3, 55, 0.9052063],
[4, 30, 0.9524169]]

```

```

In [135]: 1 import matplotlib.pyplot as plt
2 x = [tup[0] for tup in index_high_corr_list]
3 y = [tup[2] for tup in index_high_corr_list]
4
5 plt.figure(figsize=(20,10))
6 scatter = plt.plot(x,y, '.')
7 line = plt.plot(np.unique(x), np.polyd(np.polyfit(x, y, 1))(np.unique(x)), 'r')
8 plt.title('Highest Correlation of Topics/Movie Length', fontsize=30)
9 plt.xlabel("Movie Length Rank Ordered", fontsize=20)
10 plt.ylabel("Highest Topic Correlation", fontsize=20)
11 plt.show(scatter, line)

```

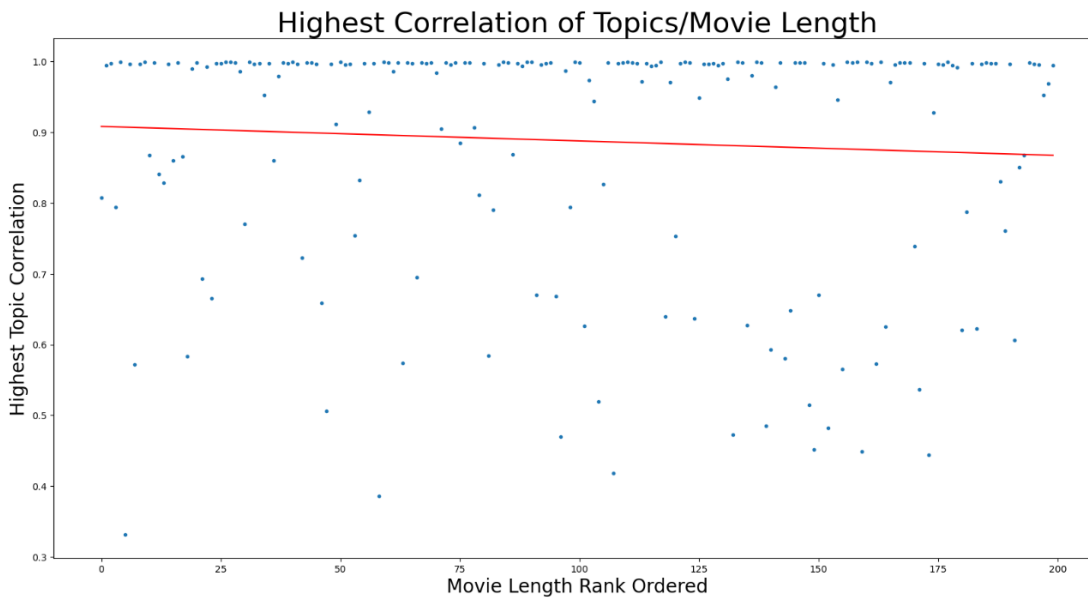


Figure 6: Βλέπουμε όντως οτι ακολουθεί πτητική πορεία

Το σύστημα συστάσεων

Είμαστε πλέον έτοιμοι να κάνουμε τις συστάσεις μας. Θα χρησιμοποιήσουμε την ομοιότητα συνημίτονου για να συγκρίνουμε τα **vectors** των θεμάτων και της συσχέτισης τους ανά έγγραφο. Η ομοιότητα συνημίτονου μεταξύ των **vectors** του σάκου λέξεων μας δίνουν ομοιότητες στην περιοχή από -1 έως 1 (όσο μεγαλύτερη, τόσο πιο παρόμοια)

$$\text{Similarity} = \cos(\theta) = (A \cdot B) / (\|A\| \|B\|), \|X\|$$

Ζητούμε αυτά τα αποτελέσματα για συστάσεις με βάση μια ταινία που έχει επιλέξει ο χρήστης. Αρχικοποιούμε τον πίνακα ομοιότητας μεταξύ όλων των εγγράφων, από τον οποίο θα ζητήσουμε για τις απαιτούμενες ταινίες και θα χρησιμοποιήσουμε ως βάση για τις συστάσεις μας. Πραγματοποιούμε τα ακόλουθα βήματα:

- Ένας τίτλος ταινίας εισάγεται από τον χρήστη
- Ο αντίστοιχος τίτλος της ταινίας βρίσκεται στον corpus
- Οι ομοιότητες όλων των άλλων ταινιών αναζητούνται από το ευρετήριο, τον πίνακα ομοιοτήτων μας
- Αυτές οι ομοιότητες απαριθμούνται έτσι ώστε μια βαθμολογία ομοιότητας να μπορεί να αντιστοιχιστεί στον αντίστοιχο τίτλο
- Οι βαθμολογίες προτάσεων για όλες τις ταινίες εισάγονται σε μια κενή λίστα που δημιουργήθηκε προηγουμένως
- Αυτές οι βαθμολογίες συστάσεων είναι ταξινομημένες
- Τα ακόλουθα εκτυπώνονται:
 - Τα πρώτα δέκα πιο συνηθισμένα **tokens** στο έγγραφο της ταινίας
 - Το σημαντικότερο θέμα της ταινίας
 - Η 2η έως την 11η συστάσεις (το πρώτο στοιχείο είναι η ίδια ταινία που ζητήθηκε)


```

In [137]: 1 def movie_recommender(title):
2     movies_checked = 0
3     for i in range(len(movies_clean)):
4         recommendation_scores = []
5         if movies_clean.iloc[i][1] == title:
6             lda_vectors = corpus_lda_model[i]
7             sims = index[lda_vectors]
8             sims = enumerate(sims)
9             for sim in sims:
10                 movies_num = sim[0]
11                 recommendation_score = [movies_clean.iloc[movies_num][1], sim[1]]
12                 recommendation_scores.append(recommendation_score)
13
14             recommendation = sorted(recommendation_scores, key=lambda x: x[1], reverse=True)
15             print("Your book's most prominent tokens are:")
16             article_tokens = corpus[i]
17             sorted_tokens = sorted(article_tokens, key=lambda x: x[1], reverse=True)
18             sorted_tokens_10 = sorted_tokens[:10]
19             for i in range(len(sorted_tokens_10)):
20                 print("Word {} (\\"{}\\") appears {} time(s)".format(sorted_tokens_10[i][0],
21                                                                     dictionary[sorted_tokens_10[i][0]],
22                                                                     sorted_tokens_10[i][1]))
23
24             print('-----')
25             print("Your book's most prominent topic is:")
26             print(lda_model_final.print_topic(max(lda_vectors, key=lambda item: item[1])[0]))
27             print('-----')
28             print('Here are your recommendations for "{}":'.format(title))
29             display(recommendation[1:11])
30
31         else:
32             movies_checked += 1
33
34     if movies_checked == len(movies_clean):
35         movie_suggestions = []
36         print('Sorry, but it looks like "{}" is not available.'.format(title))
37         other_books = []

```

```

In [138]: 1 movie_recommender("Pompeii")

Your book's most prominent tokens are:
Word 1717 ("jim") appears 90 time(s).
Word 1805 ("gang") appears 33 time(s).
Word 332 ("s") appears 30 time(s).
Word 247 ("member") appears 25 time(s).
Word 6301 ("wayn") appears 25 time(s).
Word 143 ("find") appears 18 time(s).
Word 1072 ("student") appears 17 time(s).
Word 162 ("get") appears 16 time(s).
Word 29 ("back") appears 15 time(s).
Word 1848 ("train") appears 15 time(s).
-----
Your book's most prominent topic is:
0.029*"jim" + 0.020*"s" + 0.012*"gang" + 0.010*"simon" + 0.008*"member" + 0.008*"find" + 0.008*"wayn" + 0.007*"mcclane" +
0.007*"get" + 0.007*"kate"
-----
Here are your recommendations for "Pompeii":

[['Red Riding Hood', 0.99961674],
 ['Charly', 0.99961674],
 ['Constantine', 0.99961674],
 ['13 Seconds', 0.99961674],
 ['Eye in the Sky', 0.7193544],
 ['Three Days of the Condor', 0.17535248],
 ['The Cable Guy', 0.02768309],
 ['Kolja', 0.02768309],
 ['Kingdom Hearts: Chain of Memories', 0.02768309],
 ['Attack of the Crab Monsters', 0.02768309]]

```

Figure 7: Τα αποτελέσματα της αναζήτησης

Το σύστημα συγκρίσεων

Ολοκληρώνοντας την εφαρμογή μας προσθέτουμε μια τελευταία συνάρτηση. Πιθανόν ο χρήστης να θέλει να μάθει πόσο όμοιες κρίνει το μοντέλο μας δυο ταινίες. Δημιουργούμε λοιπόν μια συνάρτηση η οποία δέχεται δυο ταινίες για είσοδο. Στην συνέχεια εκτυπώνει το ποσοστό ομοιότητας τους, ακολουθούμενο από τα πρώτα δέκα πιο συνηθισμένα **tokens** στο έγγραφο της κάθε ταινίας και το σημαντικότερο θέμα τους.

```

In [139]: 1 def movie_comparer(title1, title2):
2     movies_checked = 0
3     for i in range(len(movies_clean)):
4         similarity = []
5         if movies_clean.iloc[i][1] == title1:
6             lda_vectors1 = corpus_lda_model[i]
7             sims = index[lda_vectors1]
8             for j in range(len(movies_clean)):
9                 if movies_clean.iloc[j][1] == title2:
10                    lda_vectors2 = corpus_lda_model[j]
11                    similarity.append(sims[j])
12                    print('The similarity between {} and {} is: {}'.format(title1, title2, similarity[0]))
13
14                print('-----')
15                print("{}'s most prominent tokens are:".format(title1))
16                article_tokens = corpus[i]
17                sorted_tokens = sorted(article_tokens, key=lambda x: x[1], reverse=True)
18                sorted_tokens_10 = sorted_tokens[:10]
19                for i in range(len(sorted_tokens_10)):
20                    print("Word {} (\\"{}\\") appears {} time(s)".format(sorted_tokens_10[i][0],
21                                                                    dictionary[sorted_tokens_10[i][0]],
22                                                                    sorted_tokens_10[i][1]))
23
24                print('-----')
25                print("{}'s most prominent topic is:")
26                print(lda_model_final.print_topic(max(lda_vectors1, key=lambda item: item[1])[0]))
27
28                print('-----')
29                print("{}'s most prominent tokens are:".format(title2))
30                article_tokens = corpus[j]
31                sorted_tokens = sorted(article_tokens, key=lambda x: x[1], reverse=True)
32                sorted_tokens_10 = sorted_tokens[:10]
33                for i in range(len(sorted_tokens_10)):
34                    print("Word {} (\\"{}\\") appears {} time(s)".format(sorted_tokens_10[i][0],
35                                                                    dictionary[sorted_tokens_10[i][0]],
36                                                                    sorted_tokens_10[i][1]))
37
38                print('-----')
39                print("{}'s most prominent topic is:")
40                print(lda_model_final.print_topic(max(lda_vectors2, key=lambda item: item[1])[0]))

```

```

In [156]: 1 movie_comparer("Eye in the Sky", "Pompeii")

The similarity between Eye in the Sky and Pompeii is: 0.719353437423706

-----
Eye in the Sky's most prominent tokens are:
Word 1148 ("alejandro") appears 23 time(s).
Word 1964 ("kate") appears 23 time(s).
Word 332 ("s") appears 15 time(s).
Word 1078 ("team") appears 14 time(s).
Word 390 ("tell") appears 10 time(s).
Word 9599 ("diaz") appears 10 time(s).
Word 1341 ("gun") appears 9 time(s).
Word 6057 ("silvio") appears 9 time(s).
Word 2308 ("mexican") appears 8 time(s).
Word 6016 ("fausto") appears 8 time(s).
-----
{}'s most prominent topic is:
0.029*"jim" + 0.020*"s" + 0.012*"gang" + 0.010*"simon" + 0.008*"member" + 0.008*"find" + 0.008*"wayn" + 0.007*"mcclane" + 0.007*"get" +
0.007*"kate"
-----
Pompeii's most prominent tokens are:
Word 1717 ("jim") appears 90 time(s).
Word 1805 ("gang") appears 33 time(s).
Word 332 ("s") appears 30 time(s).
Word 247 ("member") appears 25 time(s).
Word 6301 ("wayn") appears 25 time(s).
Word 143 ("find") appears 18 time(s).
Word 1072 ("student") appears 17 time(s).
Word 162 ("get") appears 16 time(s).
Word 29 ("back") appears 15 time(s).
Word 1848 ("train") appears 15 time(s).
-----
{}'s most prominent topic is:
0.029*"jim" + 0.020*"s" + 0.012*"gang" + 0.010*"simon" + 0.008*"member" + 0.008*"find" + 0.008*"wayn" + 0.007*"mcclane" + 0.007*"get" +
0.007*"kate"

```

Figure 8: Τα αποτελέσματα της σύγκρισης

Τα συμπεράσματα

Ξεκινώντας την έρευνα για να βρούμε τον τρόπο υλοποίησης της εφαρμογής αυτής, στόχος ήταν επιπλέον να χρησιμοποιηθούν εξτρά δεδομένα όπως σκηνοθέτες, είδος, παραγωγοί, έτος κυκλοφορίας, καστ κ.α. Δυστυχώς η φύση του αλγορίθμου είναι τέτοια που βασίζεται στο πλήθος και την επανάληψη των λέξεων για να αποδώσει σωστά. Από τα αποτελέσματα που λαμβάνουμε συμπεραίνουμε ότι ο αλγόριθμος αποδίδει ικανοποιητικά στην κατηγορία των ταινιών. Λόγω του μικρού εύρους της σύνοψης κάθε ταινίας, δεν αποδίδει όσο καλά αποδίδει ο αλγόριθμος σε βιβλία και βιβλιογραφικές εργασίες.

Συνεπώς σε μια μεταγενέστερη μορφή της εφαρμογής προτείνουμε να χρησιμοποιηθούν τα σενάρια κάθε ταινίας αντί για την σύνοψη. Με αυτή την αλλαγή υποθέτουμε ότι τα αποτελέσματα θα είναι βελτιωμένα διότι το κείμενο θα πλησιάζει την μορφή ενός βιβλίου. Η εργασία βασίζεται στην βιβλιογραφική εργασία **”Recommendation system based on semantic scholar mining and topic modeling on conference publications”**.