

Team Number: SWS3001 Cluster2-10

BOCKSTO MOVIES

Website: https://yee172.github.io/Books2Movies/

Group Members

Tian Fengrui, XJTU

Peng Jiawei, ZJU

Ye Qihao, SUSTech

Liu Yihong, SCU

CONTENTS

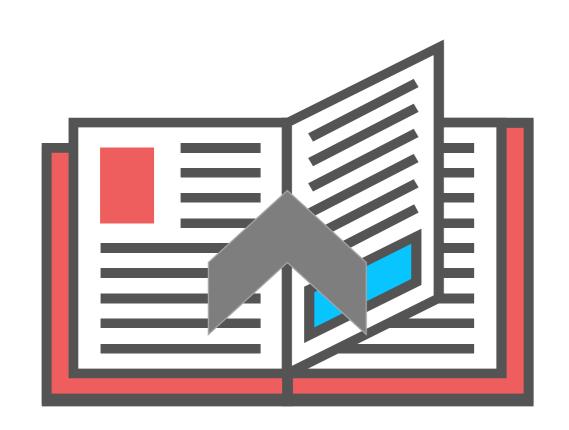
- 1 Introduction
- 2 Workflow
 - NEW Dataset
 - NEW Algorithms
 - NEW Graphs
- 3 Interesting Results
- 4 Examples

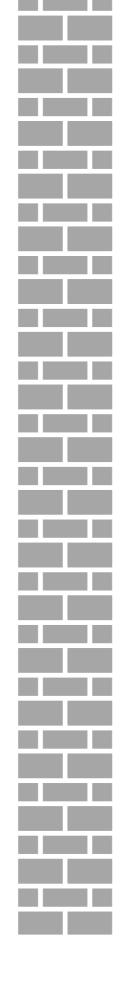
Introduction

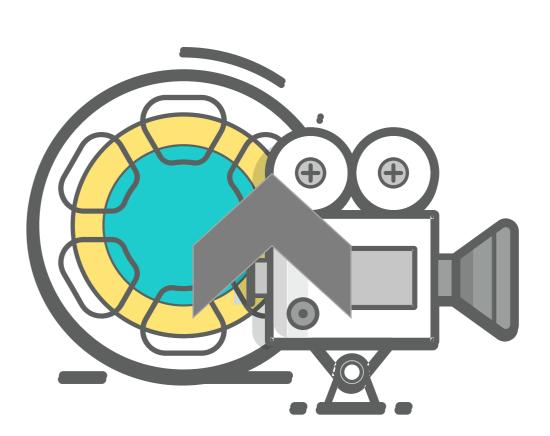
How to recommend Movies based on books

Introduction









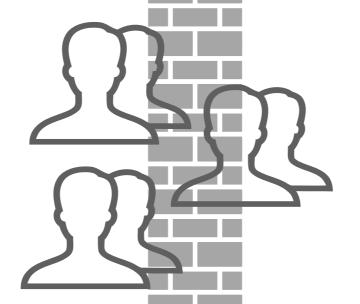
Introduction



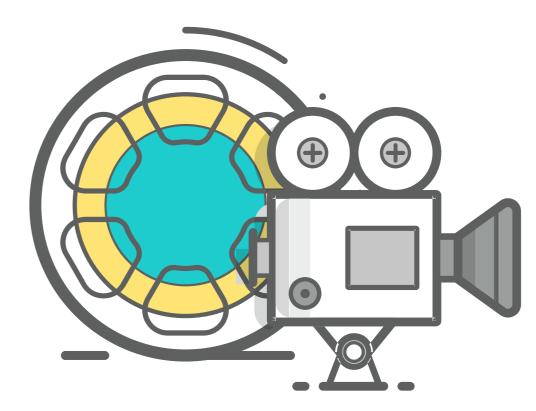
Fextre phph











Worlaflow

New dataset, new algorithms, new graphs.

Data Collecting





J Different Books
Using different tags
to separate books.
7 kinds of books in
total.

Users' Preference in Movies
Collecting each user's movies
he/she has watched.

O 2 Users' Information Collecting user' info in different book groups.

Movies' Information
Getting names, directors,
scriptwriters, actors, types,
etc.

This time to make the dataset more representative, we decide to use Chinese library classification

Big Tag	Sub Tag	Book
	诗词	苏东坡文选
文学	散文	xx散文集
	随笔	鲁迅杂文
••••		
	网络	••••
科普	编程	••••
	算法	••••

For every subtag here we find top15 books here. This may ensure the result of first graph better, and in a consequence making our findings better.



After counting out those who read all kinds of books equally frequent, then select the top 125 active users to build the graph.

B1 = Literature

B2 = Popular

B3 = Culture

B4 = Philosophy

B5 = Military

B6 = Economic

B7 = Art

B8 = Science

	User_id	Book_id	Category
4 * * * 5 * *	1	1	B1

User_id	B1	B2	••••	tag
1	1	0		Not decided



When the calculation is done, we can generate their tags, which is the top3 kinds of book they prefer.

User_id	B1	B2	B3	B4	••••	Big tag
0	11	22	1	2		B2 B1 B6
1	17	0	3	12		B7 B5 B4
0 0 0 0 0						

We generate graph1 based on this table.

National University of Singapore

School of Computing

After counting out those who read all kinds of books equally frequent, we select the top 125 active users to build the graph.

The first graph are 105 nodes and 6 clusters. We will analysis the clusters.

The number of people who watched this movie in each cluster represents how it is liked

Figure out the movies that these users have seen.



Every cluster in Graph1 may have some meanings, so we use them as the element of the movies' vector.

01		•
	L CONO	MIC
	Econo	IIIIC

C2 = Popular

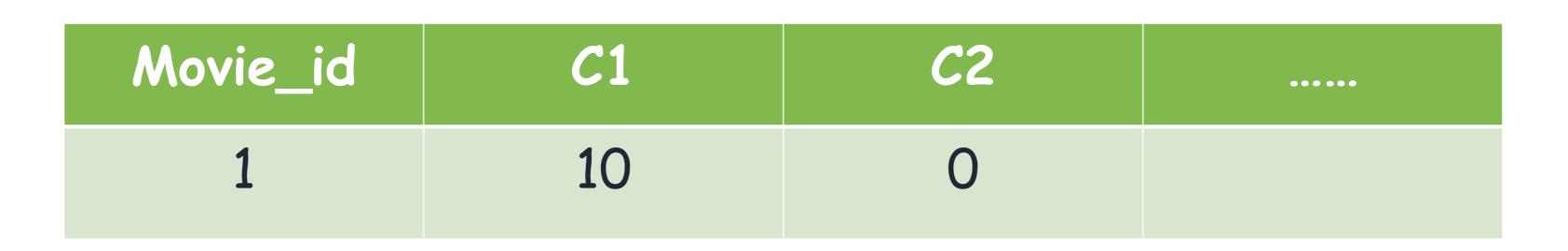
C3 = Philosophy

C4 = Science

C5 = Philosophy&Science

C6 = Literature

Cluster	Movie_id	Number of like
1	1	10





When the calculation is done, we can generate their tags according to how they are liked by the clusters, and compare them with their original tags

	Movie_id	cluster1	cluster2	cluster 3	••••	OriginTag	Cluster_t ag
Inception	0	12	7	15	••••	Science Fiction Suspicion Adventure	C 6
Life of Pi	1	3	5	12 Fantasy Adventure		C 6	
					•		

Dismat Calculation



Cosine similarity:

similarity
$$(u, v) = \frac{|\min(N(u), N(v))|}{\sqrt{|N(u)||N(v)|}}$$

Jaccard similarity:

$$similarity(u, v) = \frac{|\min(N(u), N(v))|}{|\max(N(u), N(v))|}$$

UserCF- Π F:

$$\text{similarity}(u, v) = \frac{\sum_{i \in \text{ITEM}} \frac{\min(u[i], v[i])}{\log(1 + |N(i)|)}}{\sqrt{|N(u)||N(v)|}}$$

similarity to distance:

$$\operatorname{distance}(u, v) = 1 - \operatorname{similarity}(u, v), \operatorname{similarity} \in [0, 1]$$
$$\operatorname{distance}(u, v) = \frac{1}{1 + \operatorname{similarity}(u, v)}, \operatorname{similarity} \in [0, +\infty)$$

normalization:

$$x' = \frac{x - \min(X)}{\max(X) - \min(X)}, x' \in [0, 1]$$

$$x' = \frac{x - \mu}{\max(X) - \min(X)}, x' \in (-\infty, +\infty)$$

$$x' = \log(1 + x), x' \in [0, +\infty)$$

$$x' = \log(x), x' \in (-\infty, +\infty)$$

$$x' = \frac{2 \arctan(x)}{\pi}, x' \in [0, 1)$$

$$x' = \frac{x - \mu}{\sigma}, x' \in (-\infty, +\infty)$$

$$x' = x - \mu, x' \in (-\infty, +\infty)$$

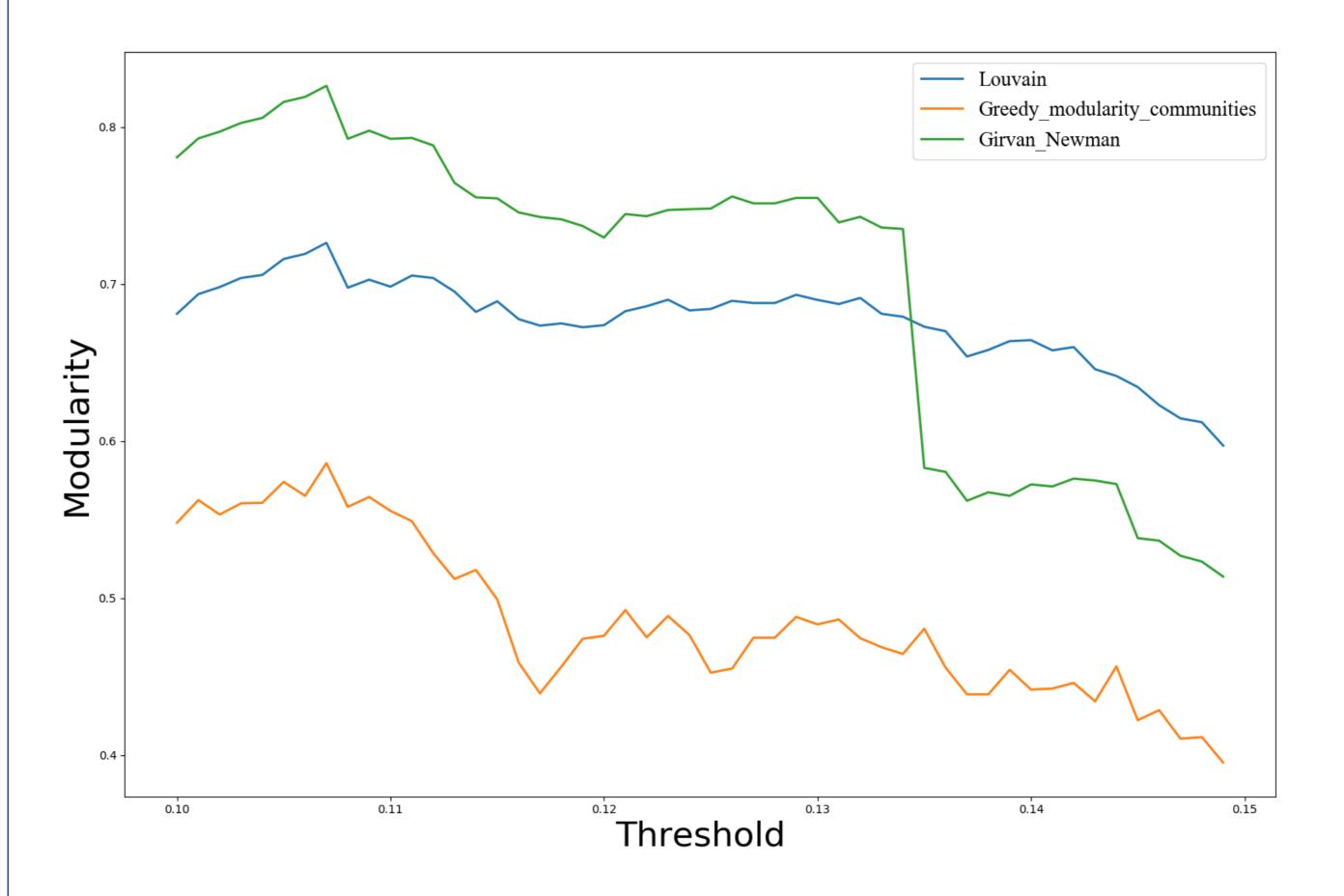
Algorithm Application



We want to choose the best algorithm and a better threshold,

so we compare the algorithm and the threshold by using this plot





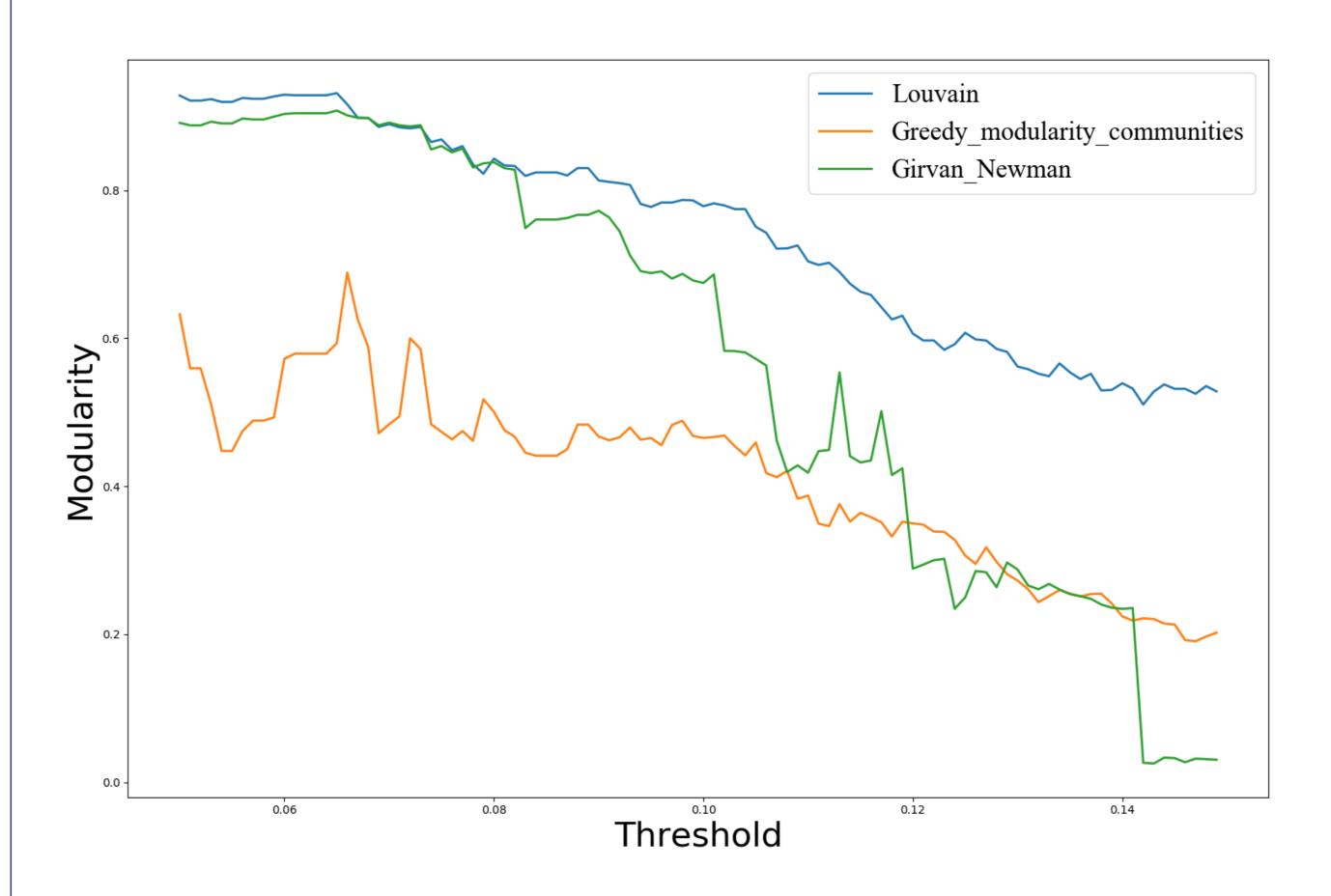
Girvan Newman works much better before the threshold goes beyond 0.135

Algorithm Application



We assume things may change when it comes to the second graph, and it really happened

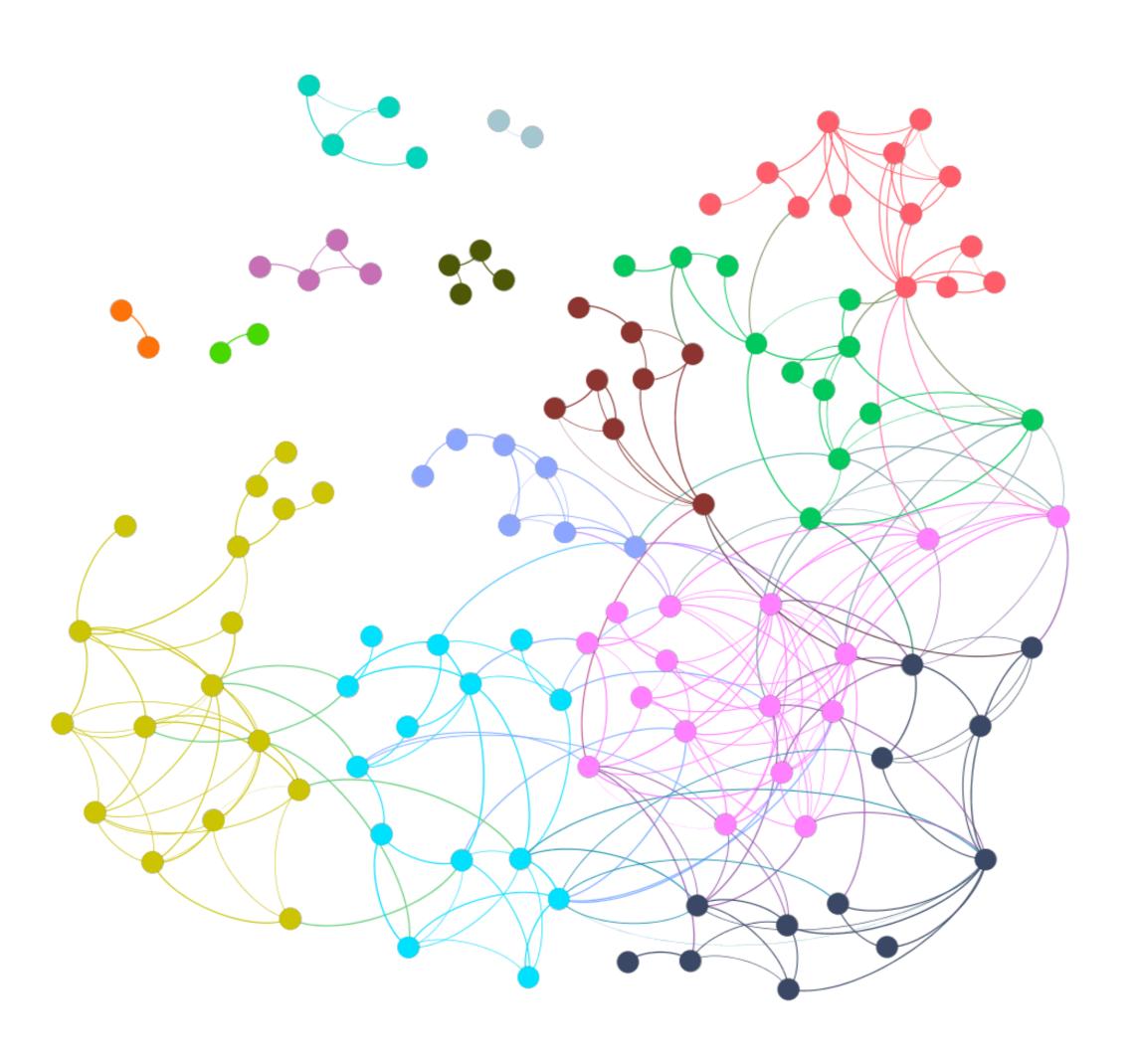




This time, Lonvain always works better







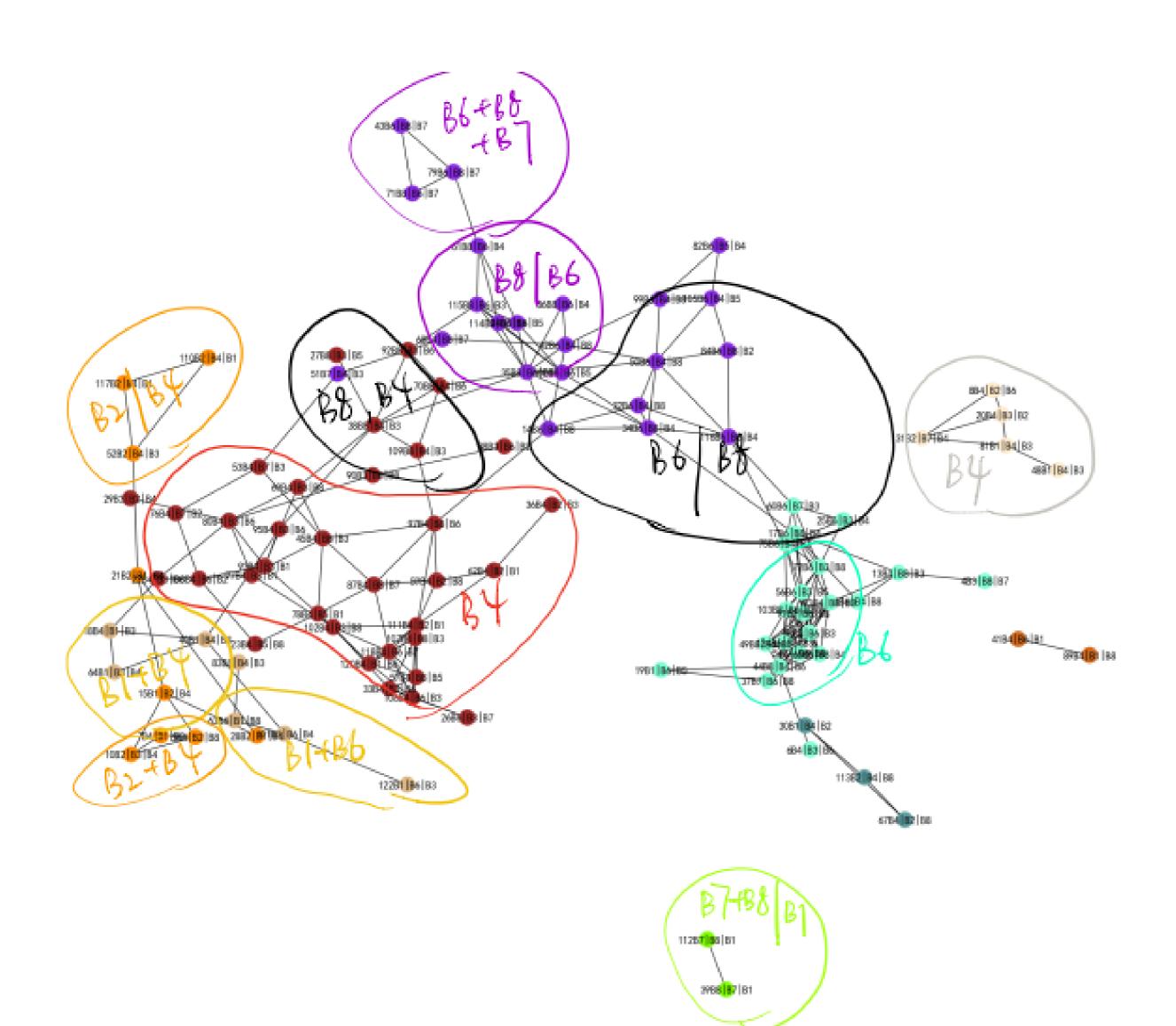
- Family + Sci-Fi
- Comedy + Romance C2+C3
- High tense but not irritating C2
- Classic
- Suspense C4
- Drama C6
- Popular
- Exciting Movie C4
- Sci-Fi C5/C1
- Popular Adventure
- Popular comedy
- Popular High rating
- Life Philosophy Sense of Time

Interesting Results

What we know we know, what we know we don't know What we don't know we know, what we don't know we don't know

Preview: Graph 1





Though the first graph only serves as a bridge from book to movie but it has some meaning.

Every community here share the same preference.
The boundary areas have a mixed preference of two kind of books.

This means the result of first graph is valid.





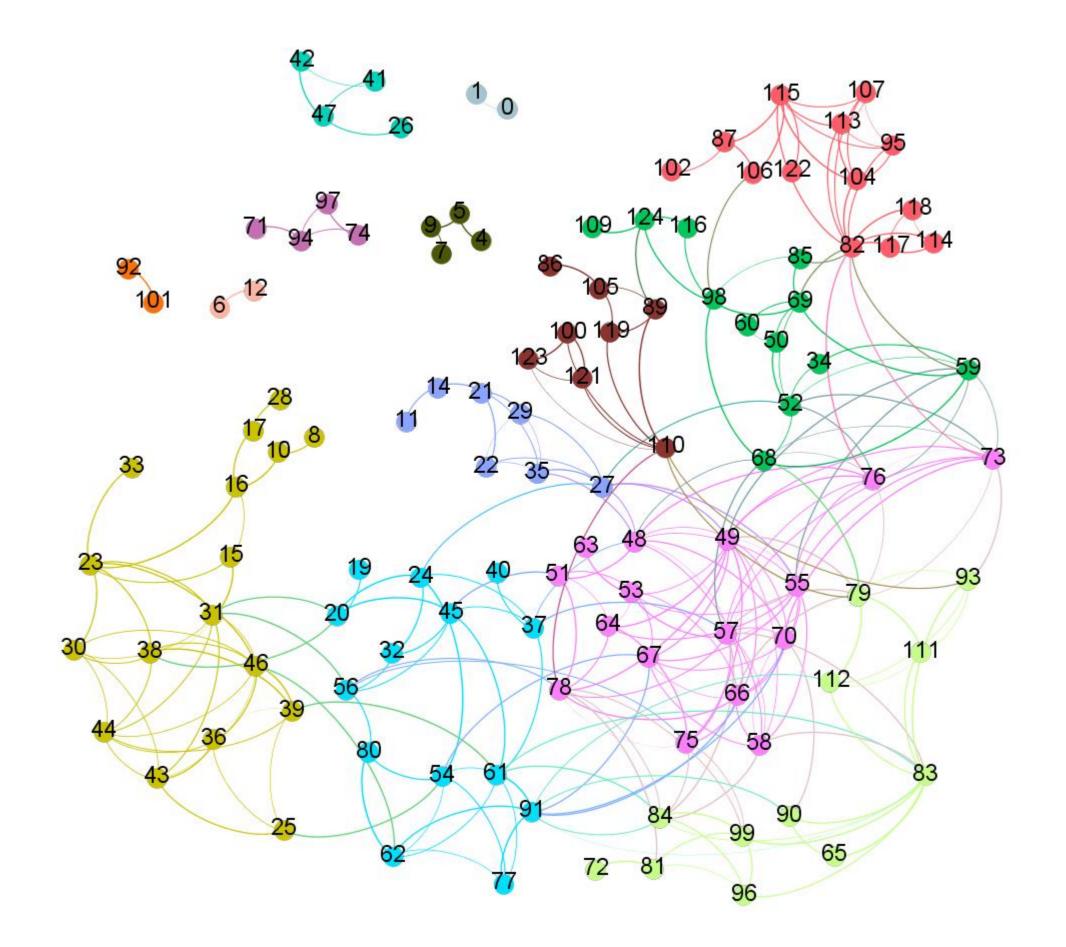
In general, many of the edges in the graph are predictable, but there are also many edges that we have not predicted. For example, the edge between *Inception* and *Life of Pi* are predictable, but the edge between *Youth* and *A Chinese Odyssey Part Two: Cinderella* is not expected.

0盗梦空间:科幻|悬疑|冒险|

29 芳华 历史 | 战争 22大话西游之大圣娶亲 喜剧 | 麦婧 | 奇幻 | 古装

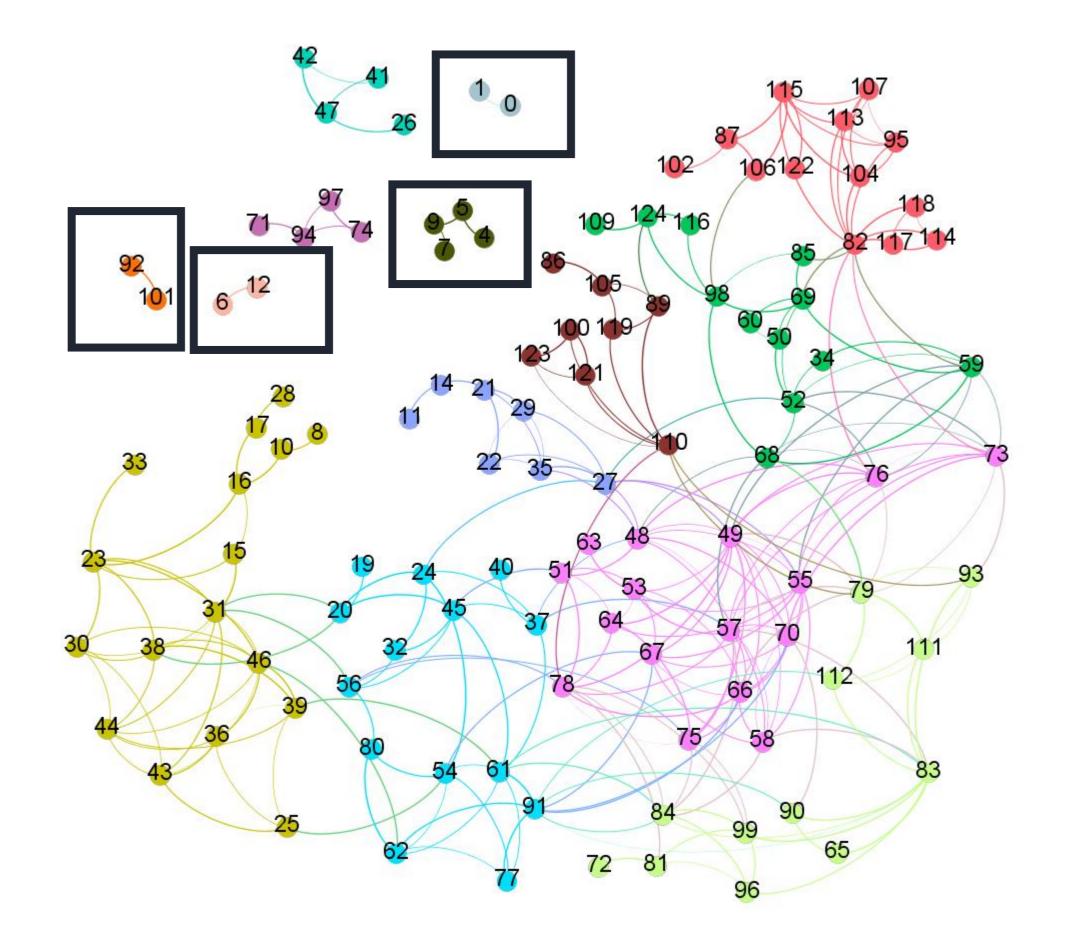


Of the most film types, there are always one or two special ones that are especially popular among people

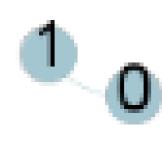




Of the most film types, there are always one or two special ones that are especially popular among people







Inception

Life of Pi

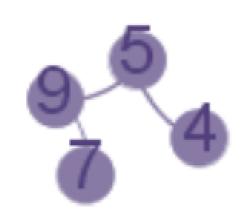
Most popular sci-fi



Let The Bullets Fly

Zootopia

Most popular comedy



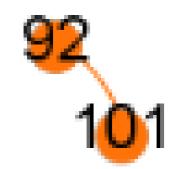
The Shawshank

Farewell My Concubine

Forrest Gump

3 idiots

TOP10 high-rating movies



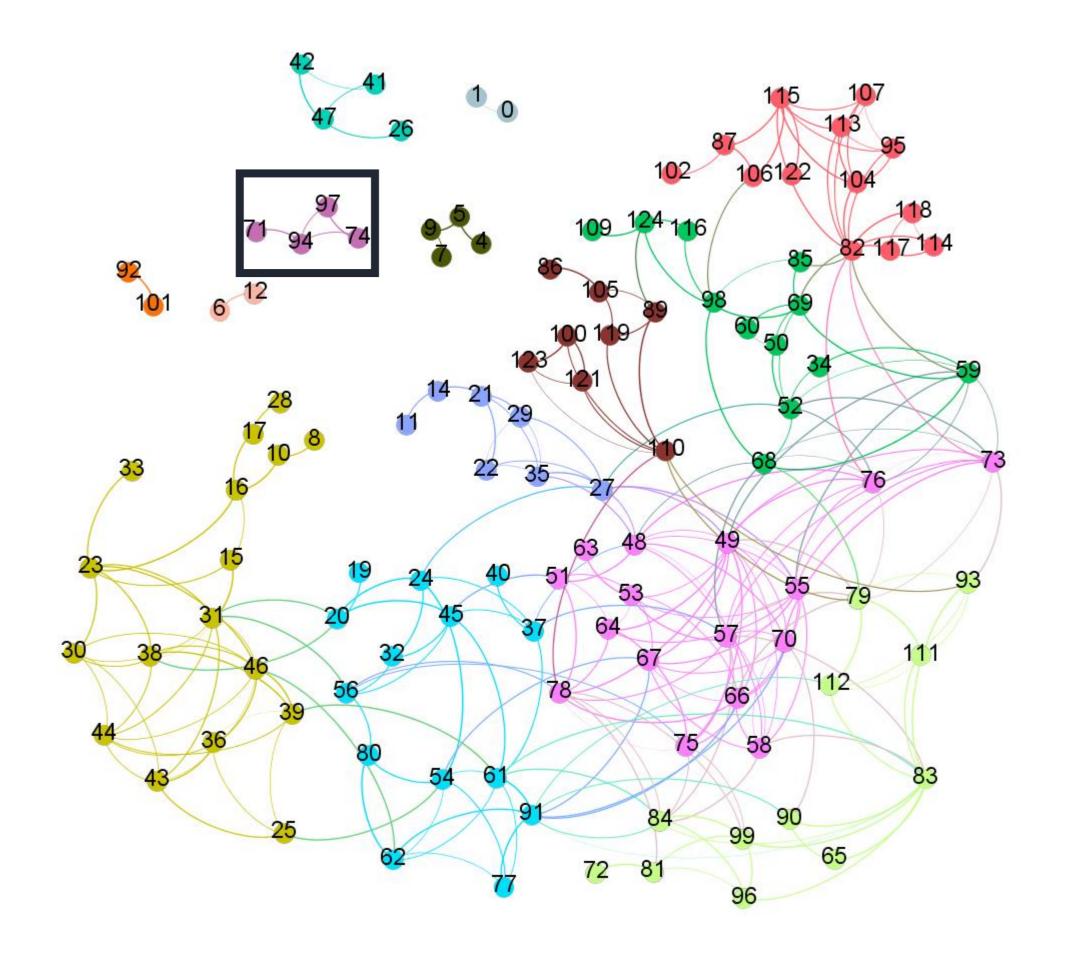
The Butterfly Effect

2012

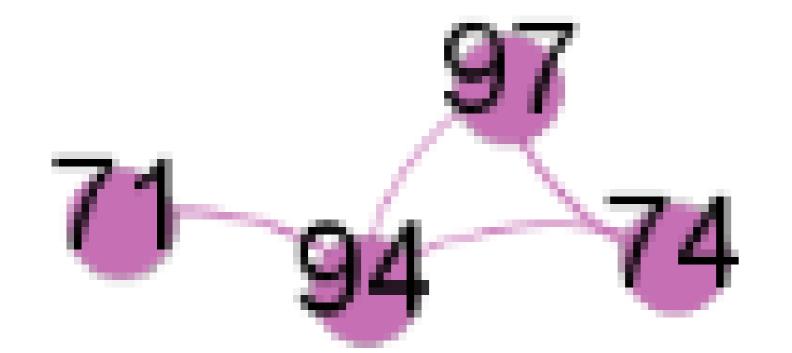
Most popular thrillers



The community that likes science books especially likes suspense.







71 Arrival

94 Avengers: Infinity

74 Train to Busan

97 Fast & Furious 7

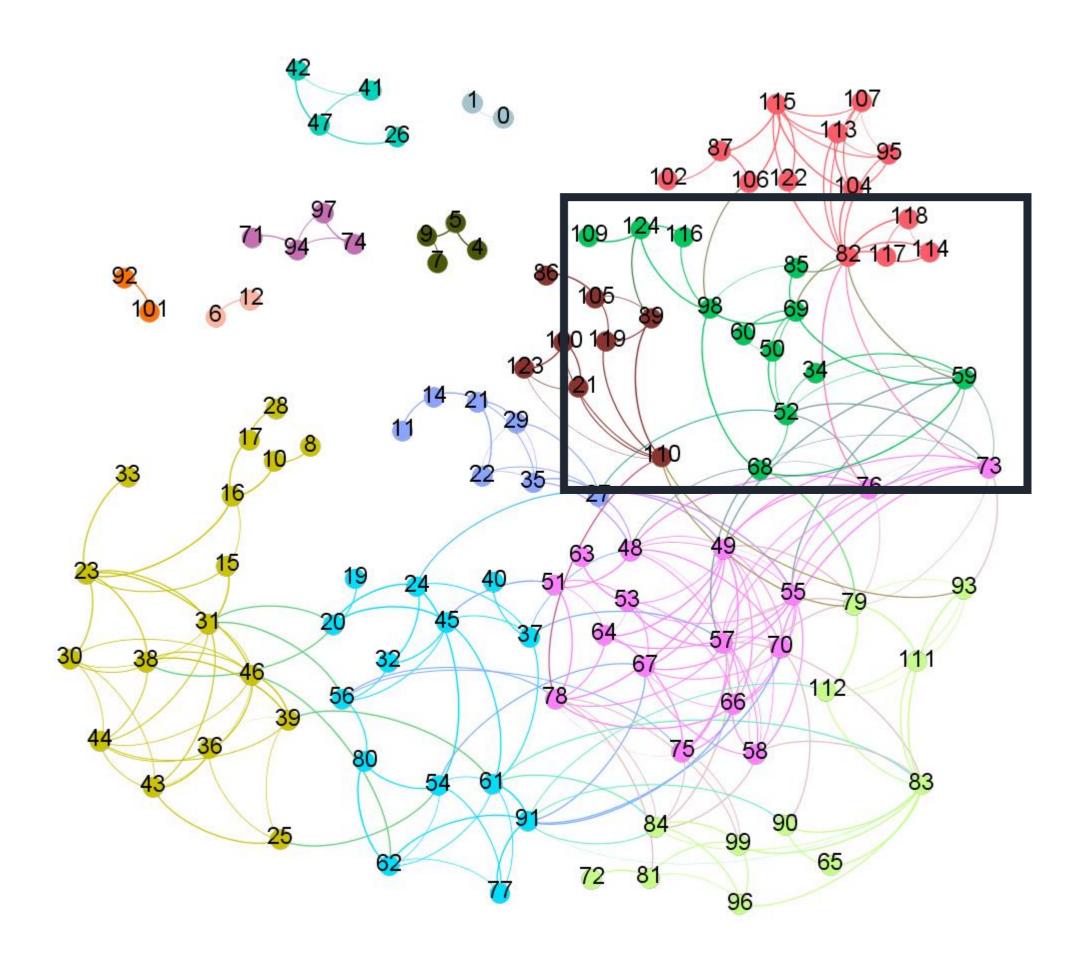
All classified because of

C4, who loves science

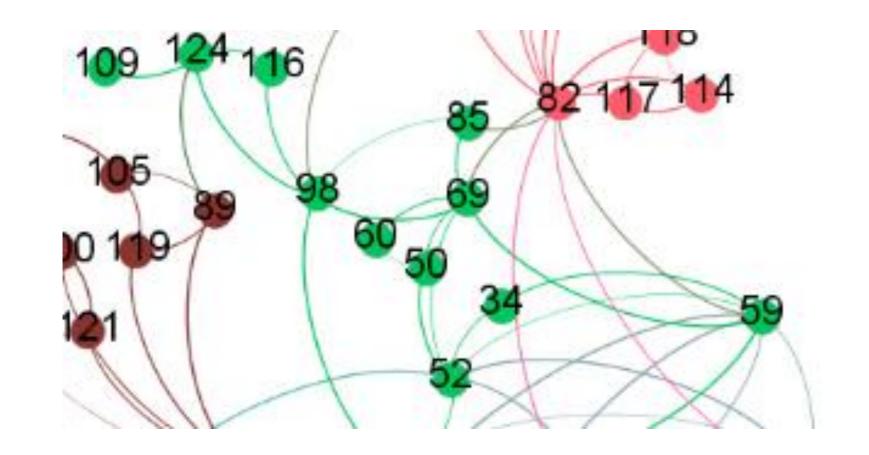
But this group is strict, they just watch the movies in the same types.



We thought those who read about pop literature may like a variety of movies, but it turns out the different.







Black Swan (黑天鹅)

Léon (这个杀手不太冷)

You Are the Apple of My Eye (那些年,我们追过的女孩)

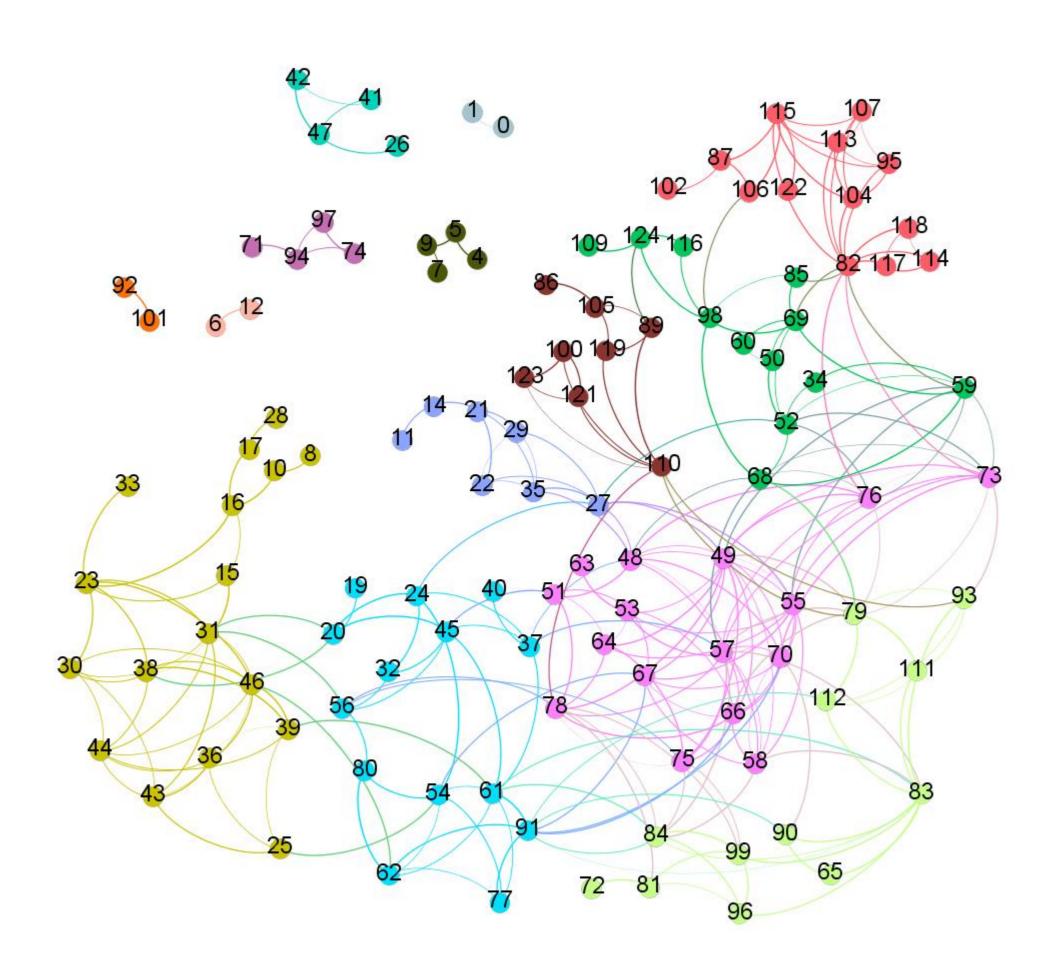
They don't like stimulating movies, but they like movies with a tight narrative, such as Black Swan, Léon, You Are the Apple of My Eye.

Actually, pop literature are mainly entertaining novels, such as novels from LouisCha(金庸) and Qiongyao, so a stimulating plot (刺激的情节) does not appeal to the readers.

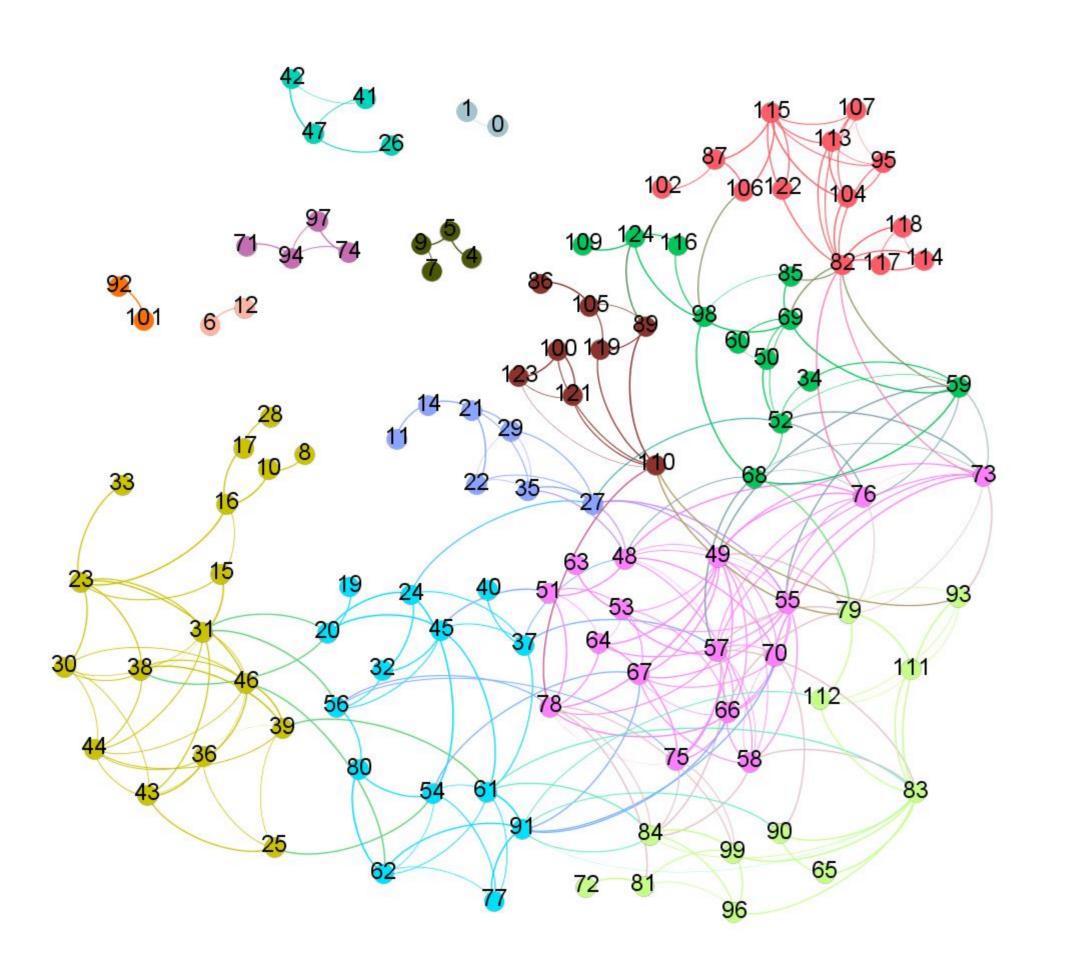
• • • • •



We thought the literary community may like only a limited number of movies. In fact, they are fond of the plot, no matter the movies' popularity, and they watch the more movies.







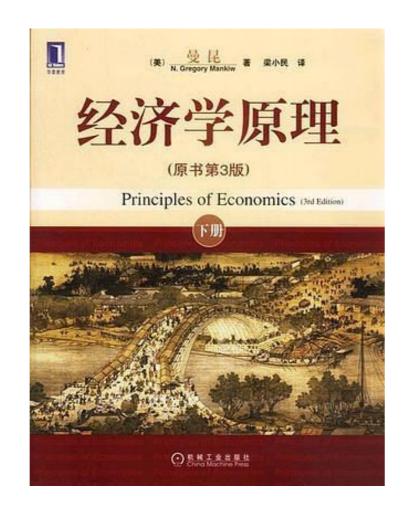
They appear in most communities.

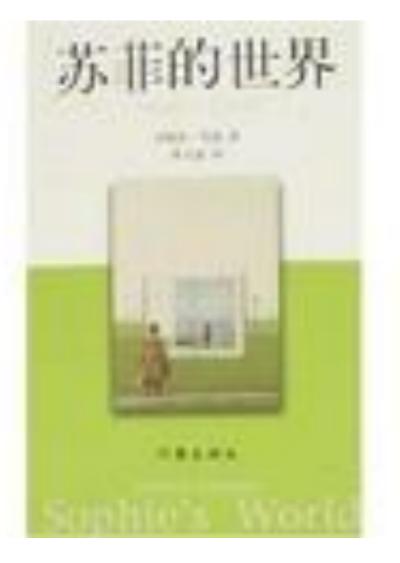
	M1	M2	M3	M4	M5	M6	M7	M8	M9	M10	M11	M12	M13	M14
C 6	0	1	0	1	1	1	1	1	1	1	0	1	1	0

So we can analysis what movie they don' like. Finally, we find out that they just don't prefer sci-fi movies. Like Avengers series and Cloud Atlas. (云图)



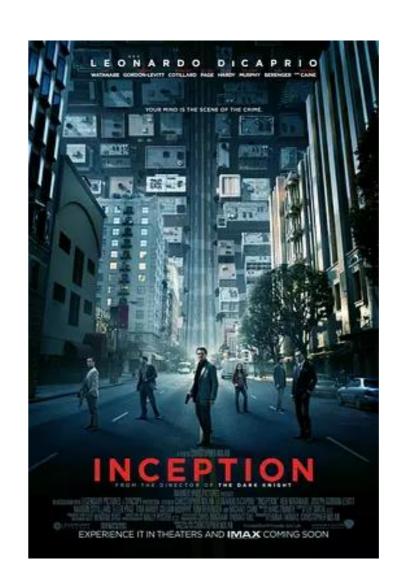
Reconsider the question: if a guy here reads these books, which movie he may prefer watching?

















The Butterfly Effect 2012

Actually those who reads about philosophy or economy like science fiction movies very much.

Even the threshold is changed ,this two node are linked together, and C1 or C5 is responsible for that.

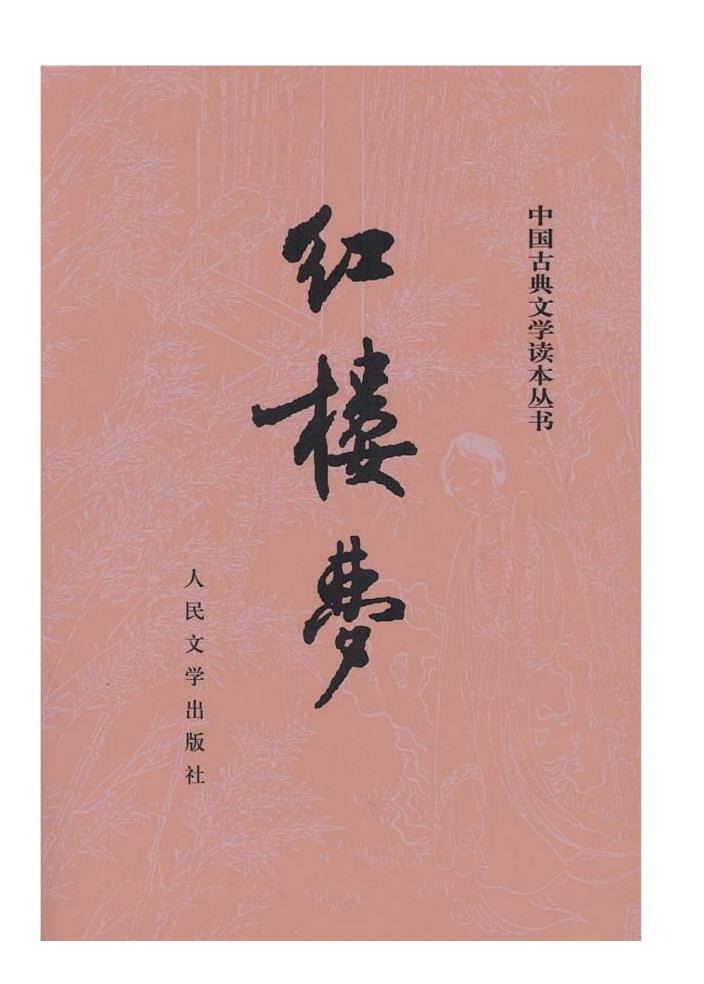
Examples

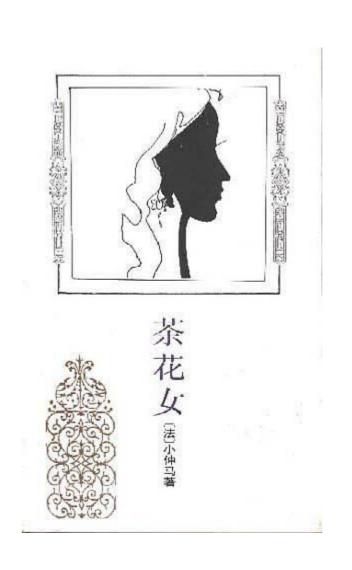
Books to Movies

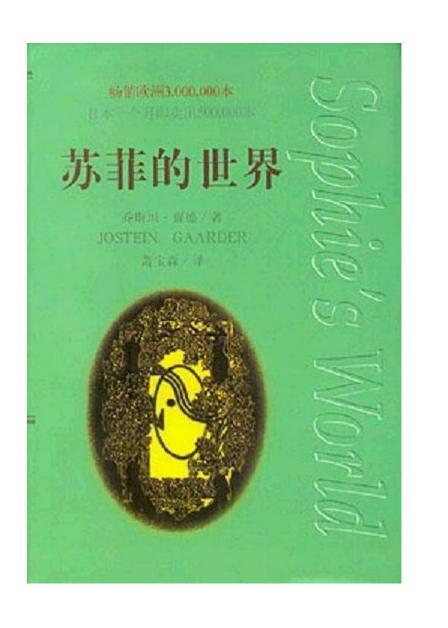






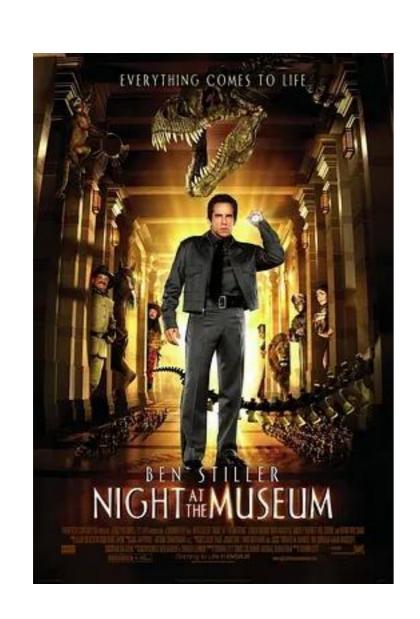












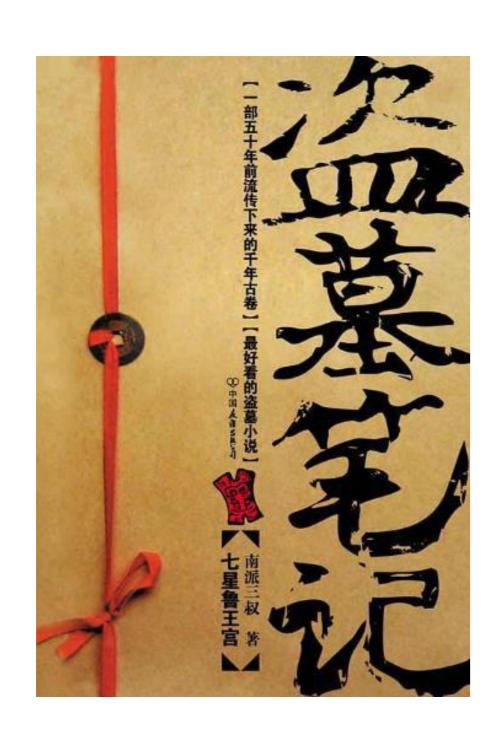




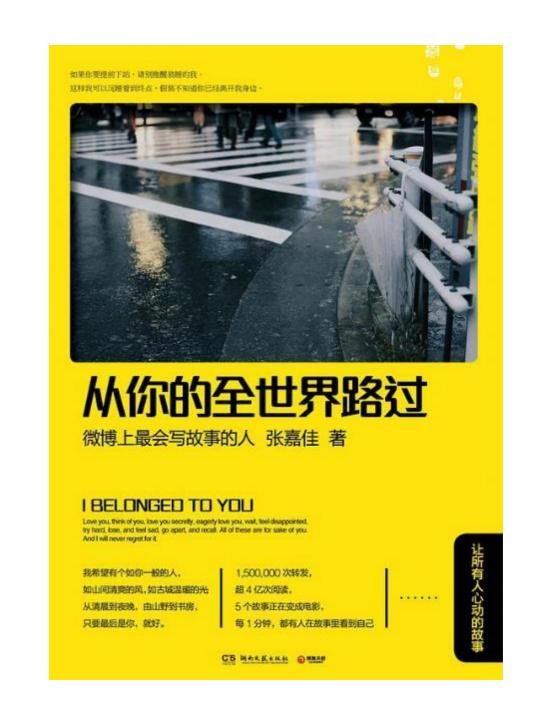










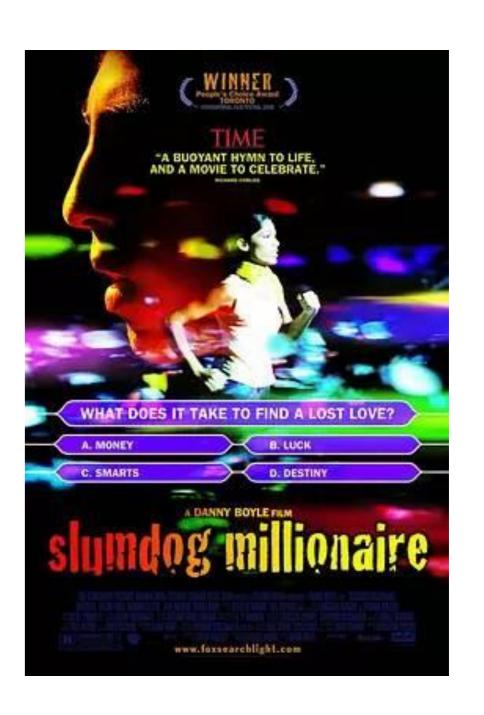






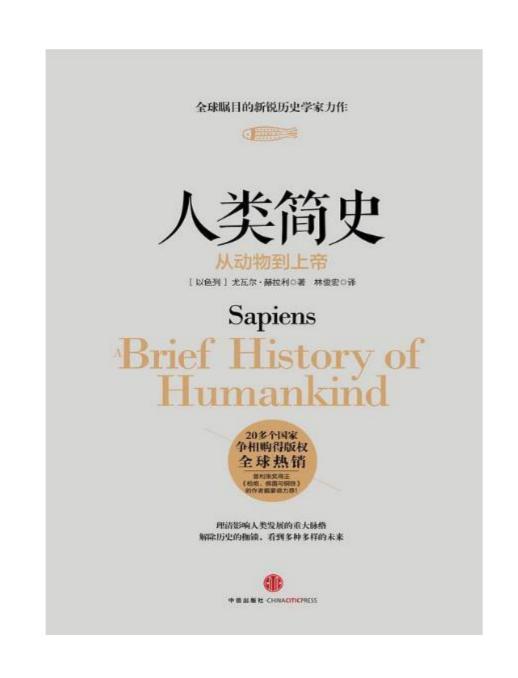




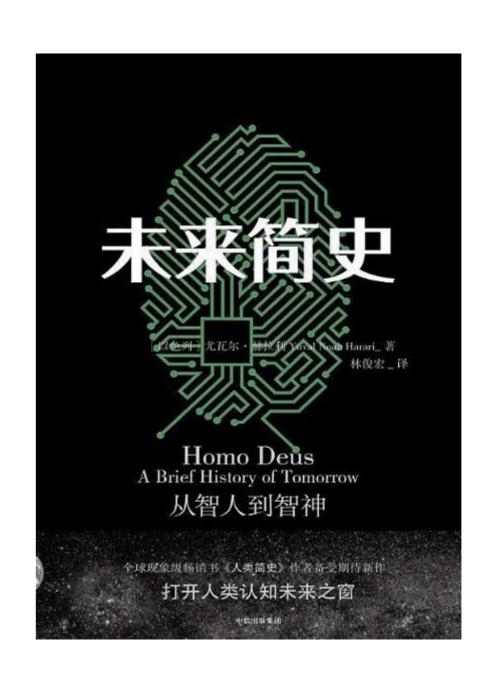








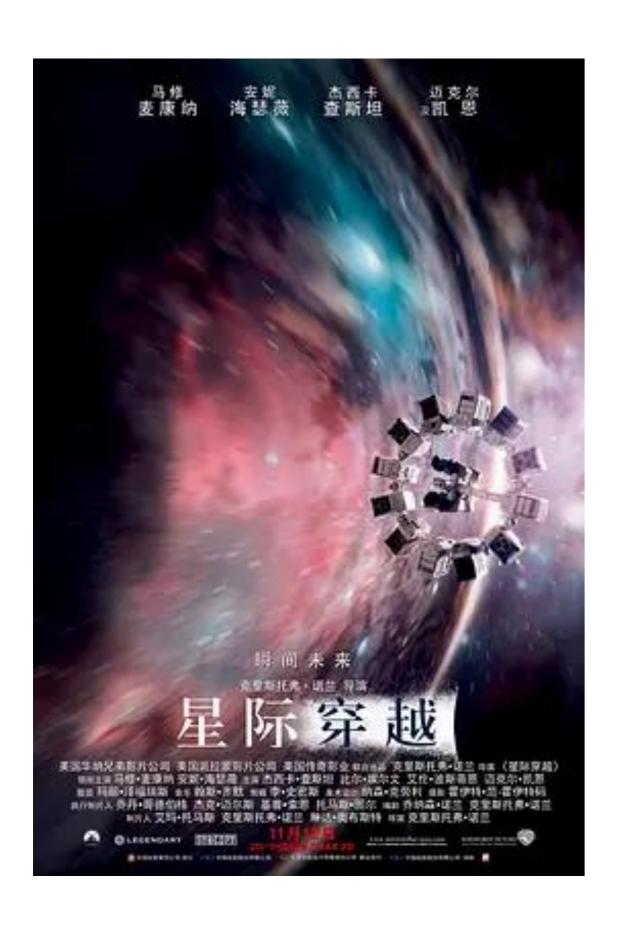




Particular about movies







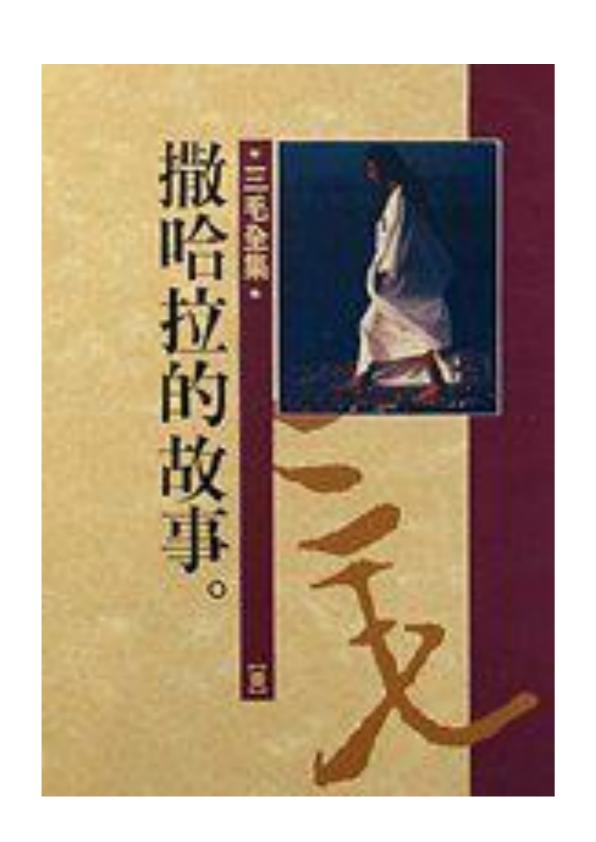




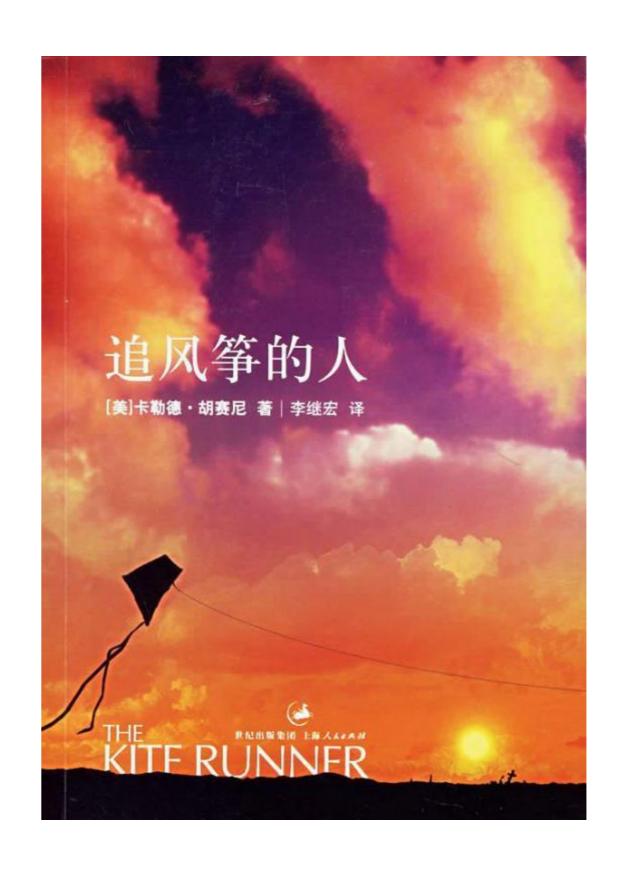
















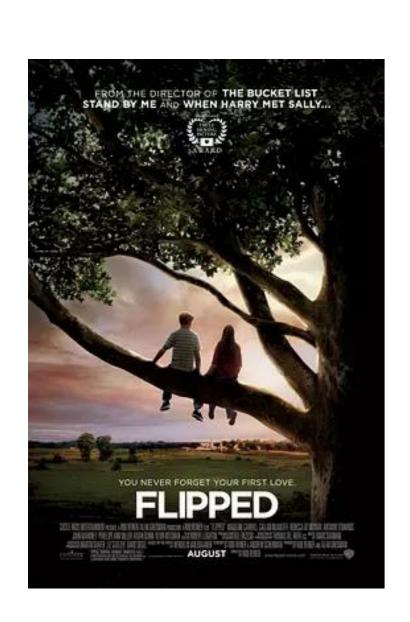


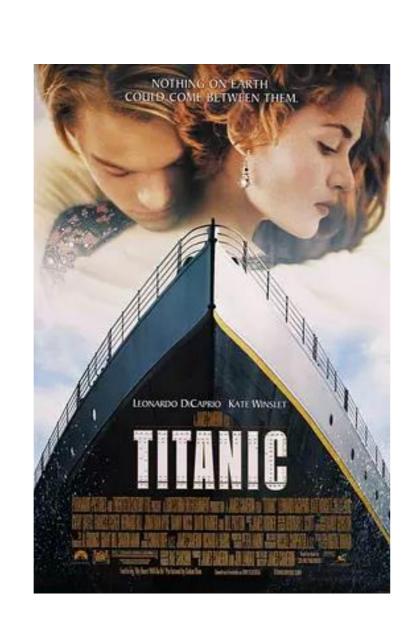




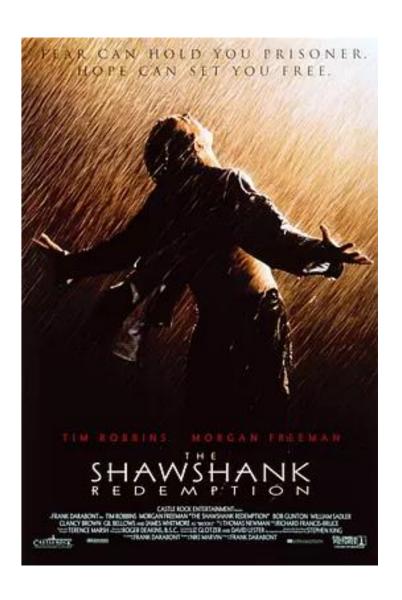














Thank you

Books to Movies