

Computational Communications and Data Analysis

Lecture 10: System Development Case Study

Ting Wang

Outlines

- Systems Thinking for Product Designing
- A Case Study: Film Box Office Prediction
- To Be A Good Data Analyst





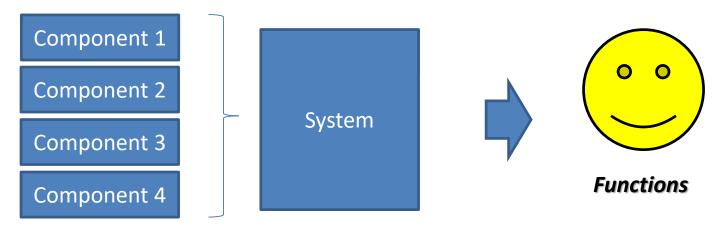


circulating development for your goals

Systems Thinking for Product Designing

What is a System?

In computer science and information science, system is a software system which has components as its structure and observable inter-process communications as its behavior.

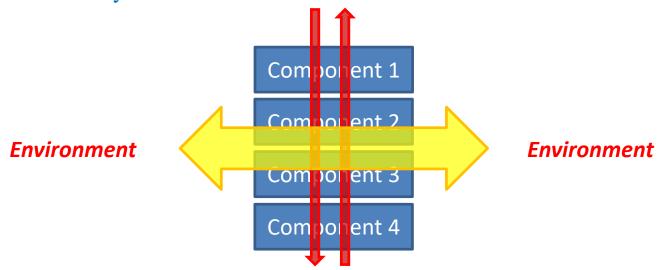




What is Systems Thinking?

Global, Optimal, and Integrated thinking methodology for software development and operation.

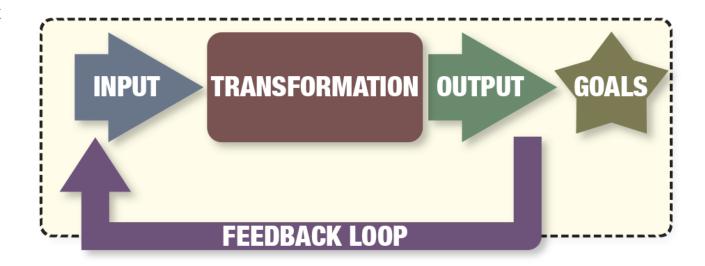
- Interactions between system and its components
- Interactions between system and its environment





Two recommended Systems Thinking Approaches

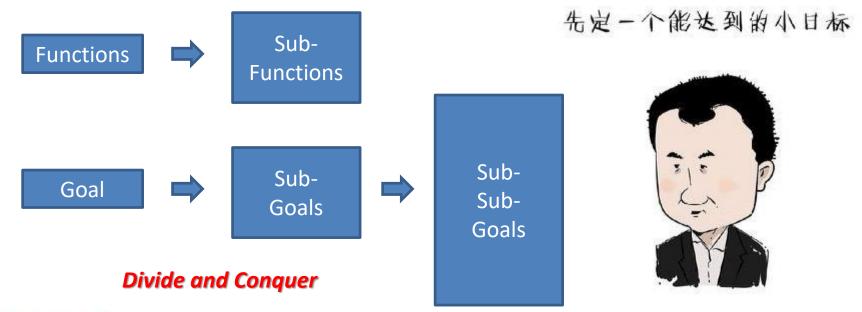
- Goal Seeking
- Input and output



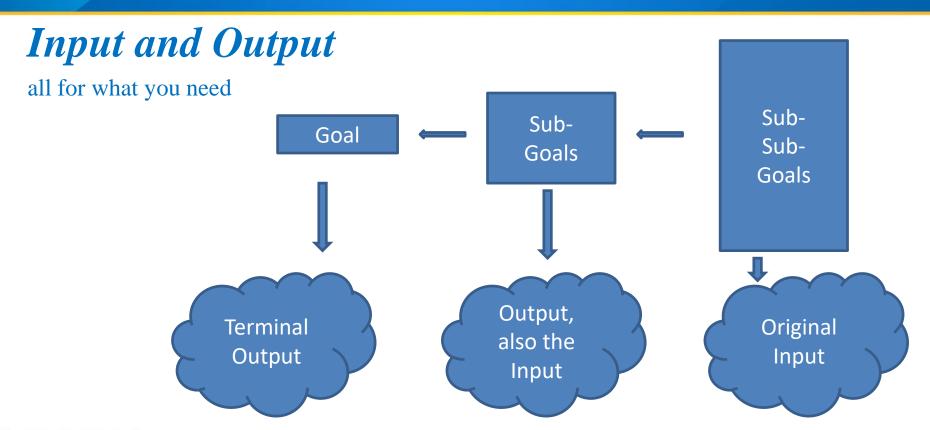


Goal Seeking (Global optimization) 全局最优

a global optimization of a function or a set of functions according to some criteria









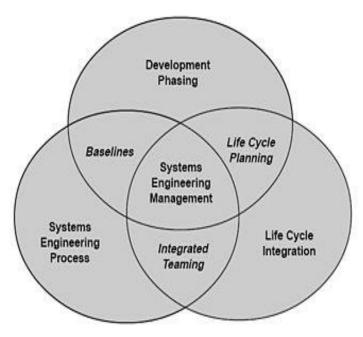
System Engineering 系统工程

ensures all likely aspects of system are considered, and integrated into a whole product.

Software Engineering

(in software and information industry)









a case study

Film Box Office Prediction



Case Description

Film Box Office Prediction

- is crucial to film investment
- is significant to the market without Completion Bond
- can be done by a number of approaches

In this case, film box office prediction will be computed based on the information collected by online film news reports.



Software Analysis



Terminal Goal

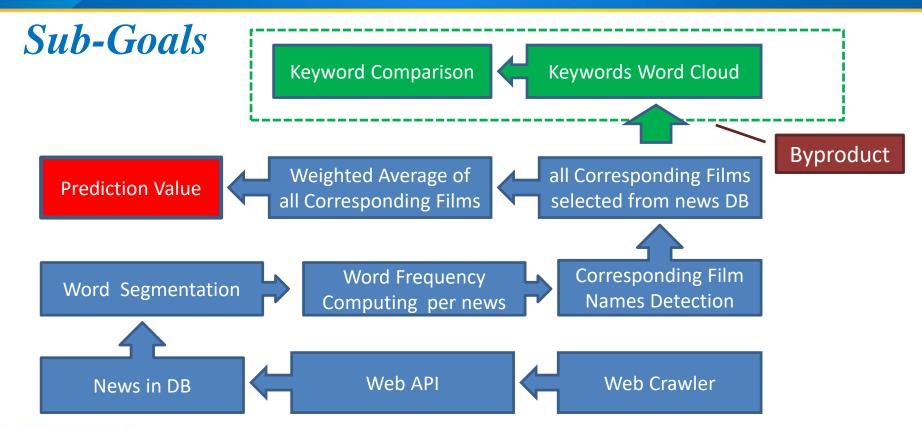
To make a decision:

whether a film is worth of being invested or not.

Final Output

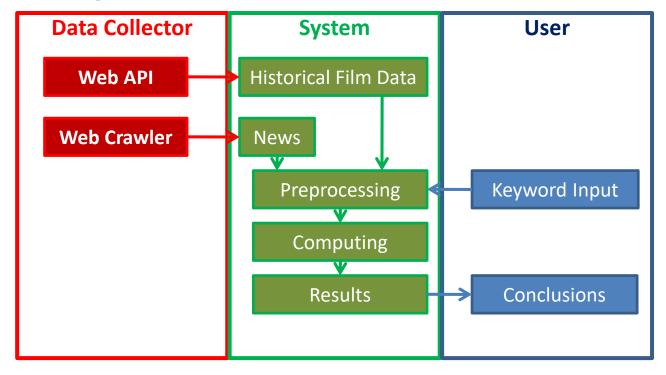
This depends on the **prediction value** of the box office of the potential film project.







Activity Diagram





Functions

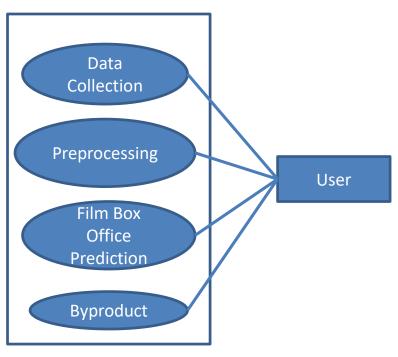
- 1. Film Box Office Prediction
- 2. Byproduct: Keyword Comparison
 - Word Cloud
 - Media Attention
 - Feature Comparisons





Use Case Diagram





Input and Output

Input: Keywords of film name

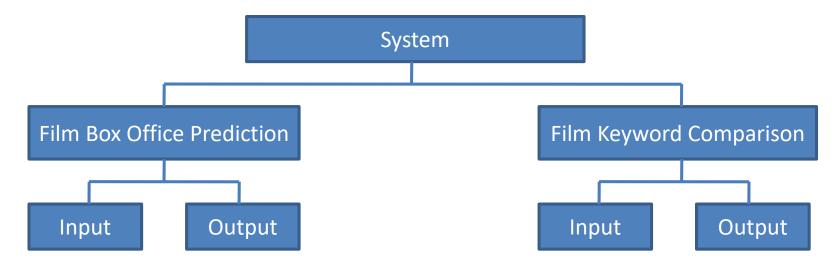
- Byproduct: Keywords
- Other conditions: Word Frequency, Periods,...

Output: Prediction value of film box office

- Word Cloud,
- Media Attention,
- Word Frequency Comparison



System Design





System Architecture

Weighted Average Computing Word Cloud, Media Attention

Historical Film Box Office Records Statistical Computing of News Report

Flask, Word Frequency Computing

Word Dictionaries

Film Box Office Prediction

Corresponding Film Detection

Keyword Feature Selection

Word Segmentation

Database

Web Crawlers

Web APIs

Preprocessing

Byproduct

News Analysis

Keyword Input

My SQL

Python



Databases

Word_Dictionary

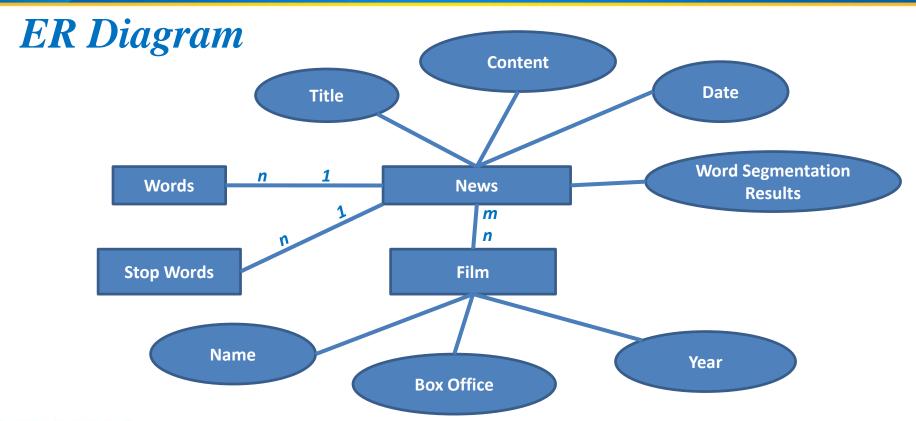
News

Stop_Word

Historical_Film_Box_Office

Tips: Film names also can be used for word segmentation.







Computing Steps

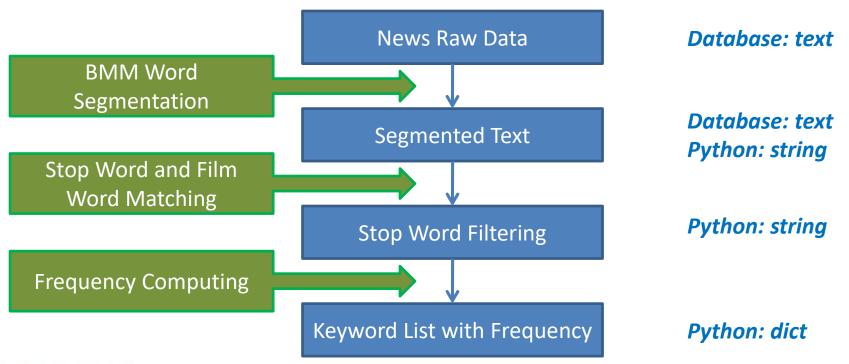


Data Collection





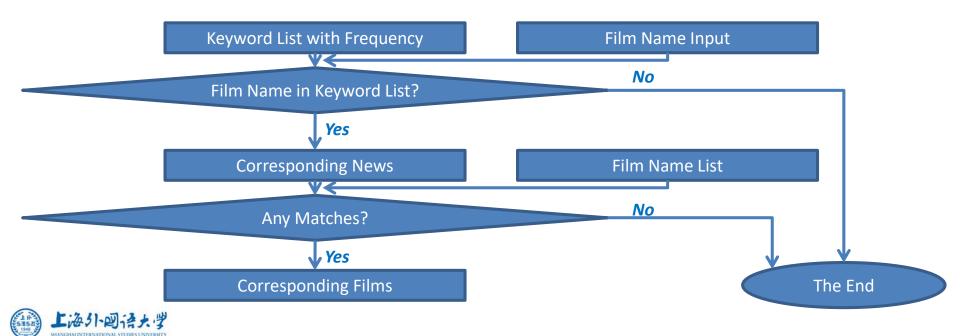
Data Transformation



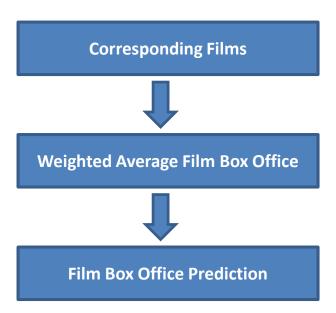


Information Acquisition (From Data to Info.)

For Film Box Office Prediction



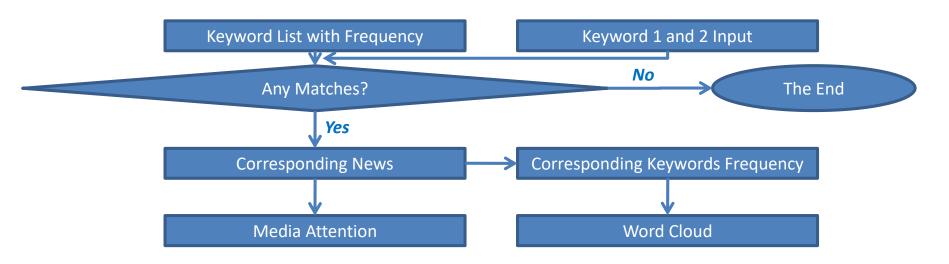
Prediction and Data Visualization



$$\overline{x} = \frac{x_1 f_1 + x_2 f_2 + \dots + x_k f_k}{n}$$

Text Mining

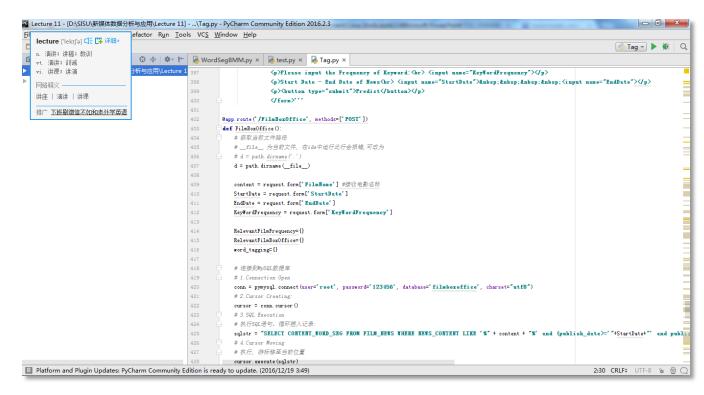
For Byproduct, Keyword Comparison





Software Development

Python PyCharm Flask MySql





Testing



Home

<u>Keyword Tagging</u>

<u>Keyword Comparison</u>

Film Box Office Prediction



Input for Keyword Comparison

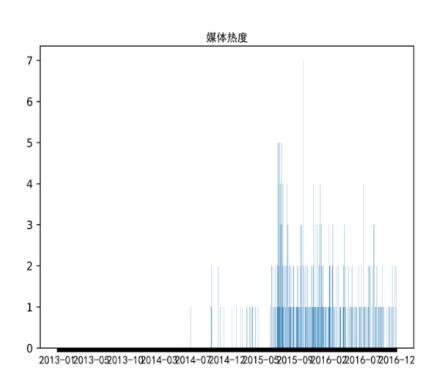
Please 捉妖记	input	the	Keywords:	西游降魔篇
wex(ra				□ 001+2 2 /m
Please	input	the	Frequency	of Keyword:
10				
Start I	Date -	End	Date	
2013-1-1	1			2016-12-1
Compa	rison			

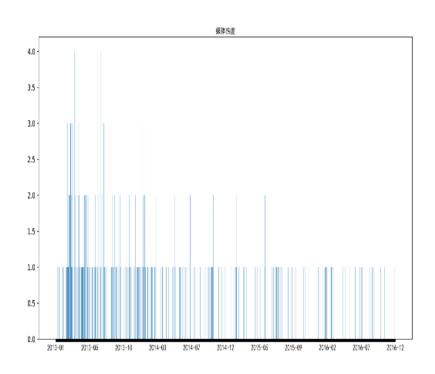


Similarity: 59.74025974025974%

Key Word 1: 芈月传 Key Word 2:甄嬛传

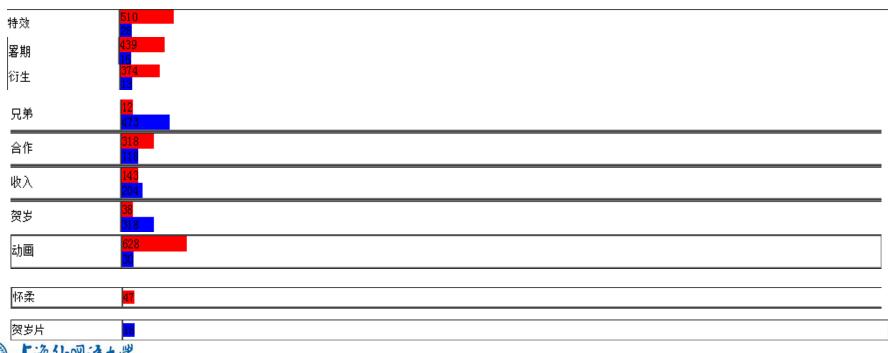








Keyword Comparison





```
← → C ① 127.0.0.1:5000/FilmBoxOffice
Please input the Film Name:
长城
Please input the Frequency of Keyword:
Start Date - End Date of News
2016-1-1
                                  2016-12-1
```

Predict





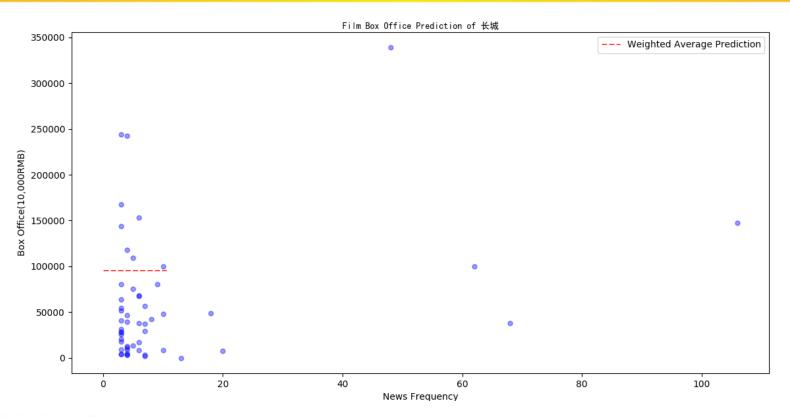
① 127.0.0.1:5000/FilmBoxOffice

Home

Film Box Office of 长城: 95428.38819320215(x10,000) RMB





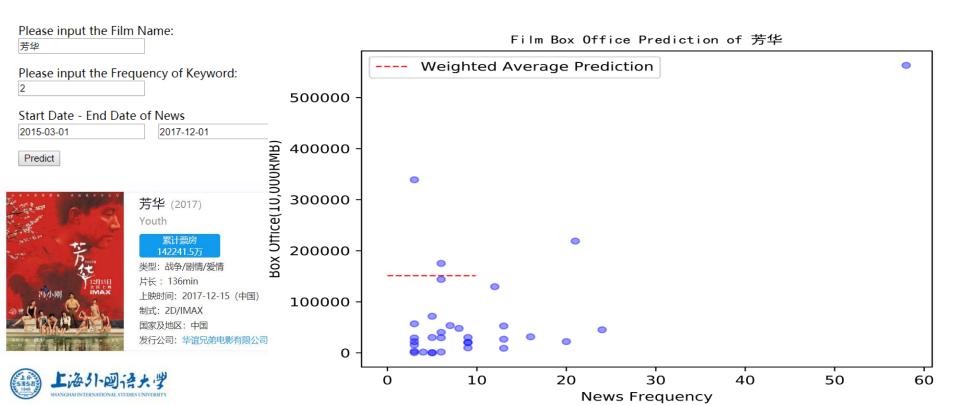




• 《芳华》

Home

Film Box Office of 芳华: 151097.2136392405(x10,000) RMB

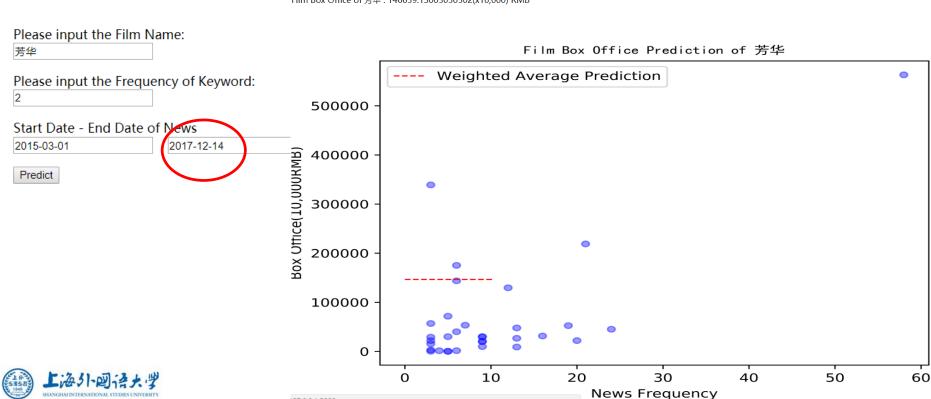




Home

127.0.0.1:5000

Film Box Office of 芳华: 146639.15003030302(x10,000) RMB



Conclusions





What are the shortages of this system?



Do you have any ideas about developing a better one?





tips for your career

To Be A Good Data Analyst

Tip 1

- You have opinions, so do data
- How to read and interpret these data is very important, it depends on your opinions
- Sometimes, GUESS is important, a hypothesis is crucial to the problem



Guess for Hypothesis



Guess for Hypothesis

哪种关系更稳定? What kind of relationship is more steady between Male and Female?

- 不是东风压倒西风,就是西风压倒东风 One Strong, One Weak
- 两种风差不多强劲 Equal

Take Films Stars as an example:

Hypothesis

男女之间,不是东风压倒西风,就是西风压倒东风,你待她 太好,她未必会投桃报李。

——司溟 《鸩ź

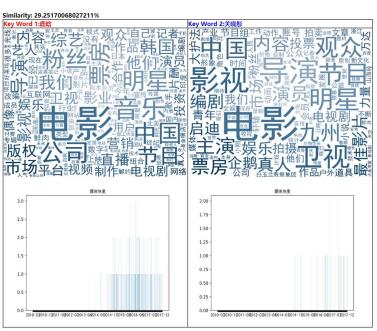


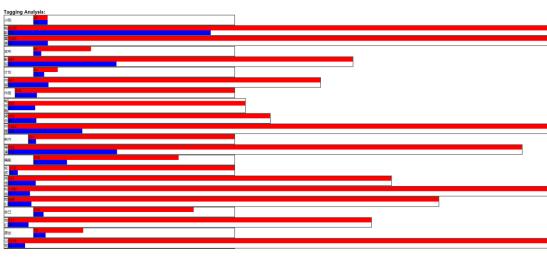
Guess for Hypothesis

- 鹿晗 关晓彤;
- 孙俪 邓超;
- 佟丽娅 陈思诚;
- 李小璐 贾乃亮



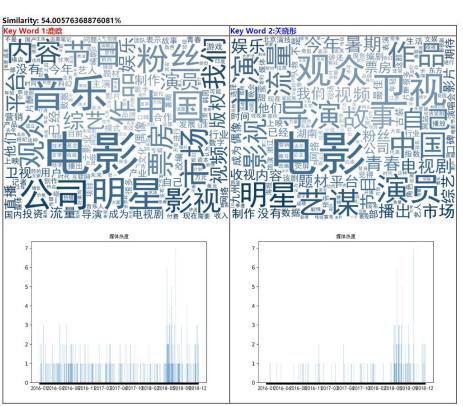
• 鹿晗 关晓彤 (2018)

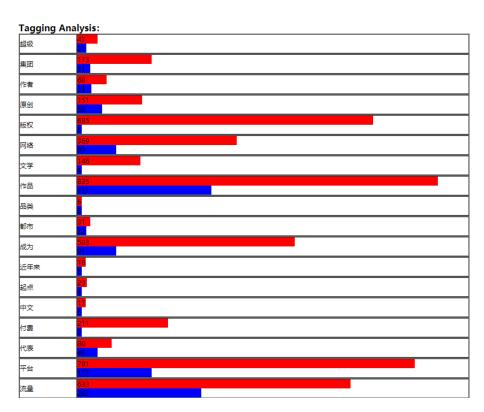




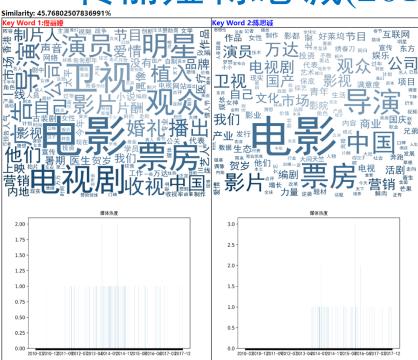


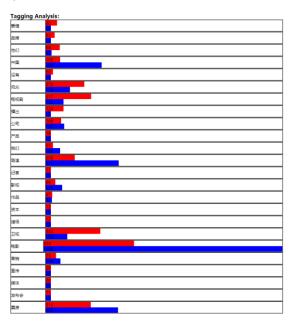
• 鹿晗 关晓彤 (2019)

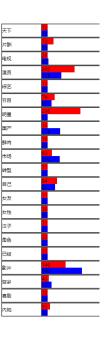




• 佟丽娅 陈思诚(2018)

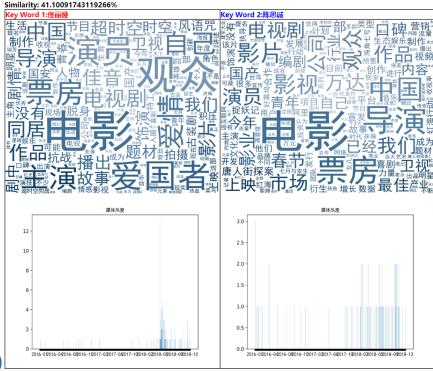


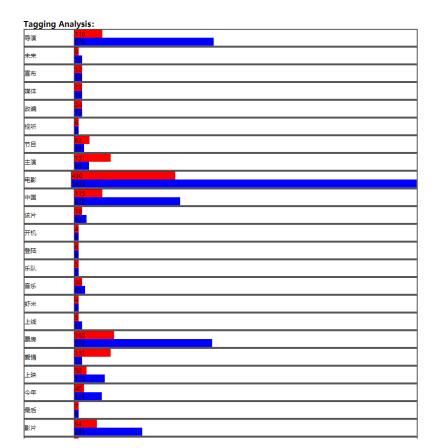






• 佟丽娅 陈思诚(2019)

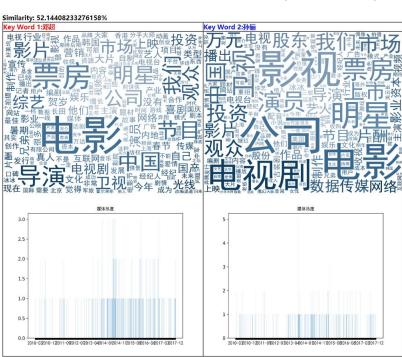


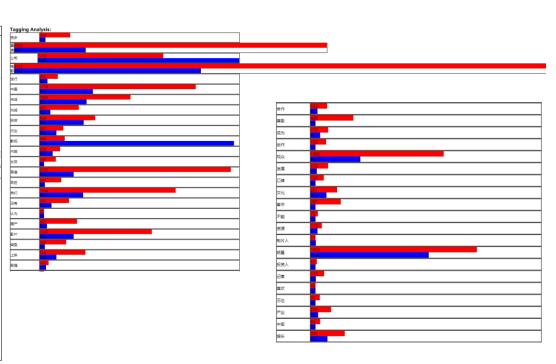




SHANGHALINTERNATIONAL STUDIES UNIVERSITY

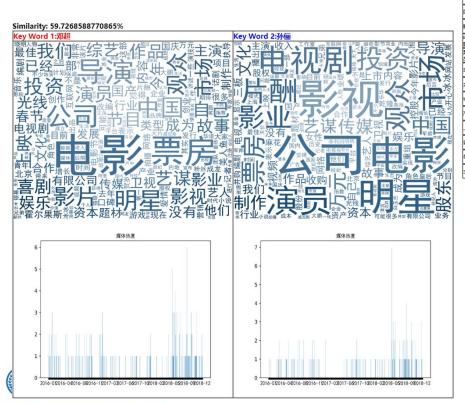
• 邓超 孙俪(2018)

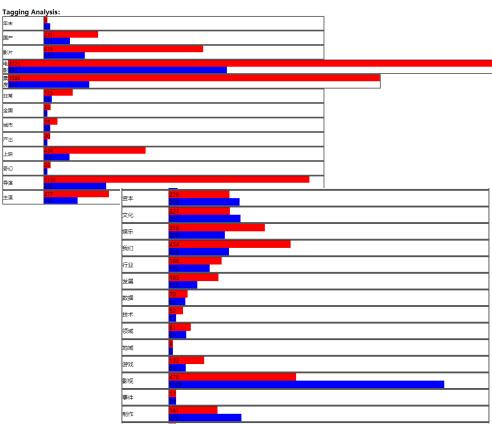




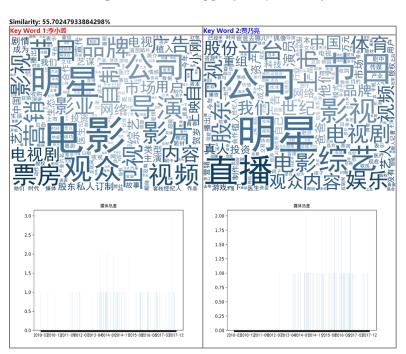


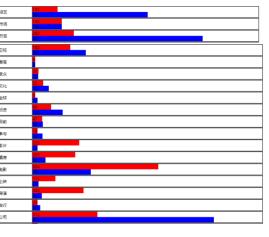
• 邓超 孙俪(2019)

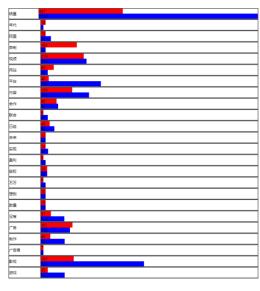




• 李小璐 贾乃亮 (2018)

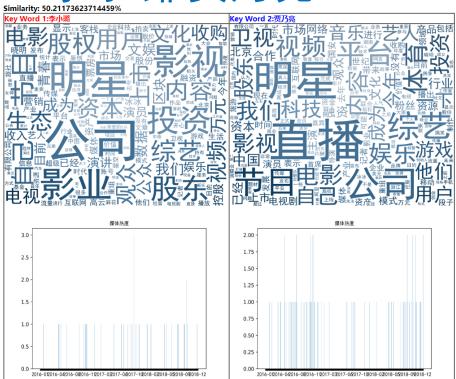


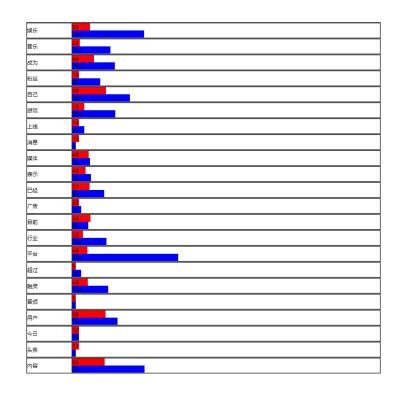






• 李小璐 贾乃亮(2019)







SHANGHAI INTERNATIONAL STUDIES UNIVERSITY





Now, what is your conclusion?



Tip 2

- Data Quality is always the most important
- Precise Prediction needs good data quality



Tip 3

• Data Analysis is not the end, but a new start. Decision Support is more important.



Tip 4

• To know more about your business, which is more important than to know more algorithms and mathematic models.



Tips 5

- Conclusions that are not correct, feasible or applicable are useless
- Conclusions will change, if some elements, such as hypothesis, time, and place are changed





Books and Chapters (1)

https://item.jd.com/11983227.html

Chapter 1-2

Machine Learning Package Installation

Machine Learning Theory Foundations





Books and Chapters (2)

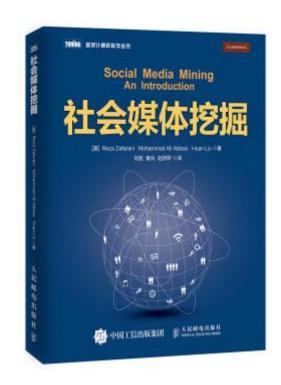
https://item.jd.com/11803260.html

Chapter 5

Data Mining Essentials

Online Reference:

http://www.public.asu.edu/~huanliu/

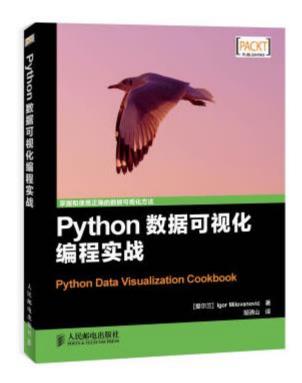




Books and Chapters (3)

https://item.jd.com/11676691.html

Python Data Visualization

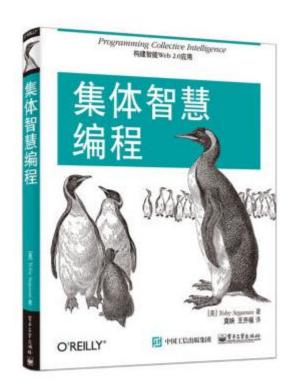




Books and Chapters (4)

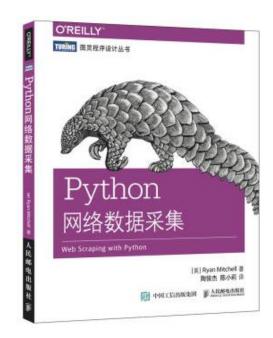
https://item.jd.com/11667512.html

Programming Collective Intelligence



Books and Chapters (5)

https://item.jd.com/11896401.html Python网络数据采集



All References for this Course:

- 张海藩.软件工程导论(第六版)[M].北京:清华大学出版社.2013年
- Meliir Page-Jones.UML面向对象设计基础[M].北京:人民邮电出版社.2012年
- 王珊、萨师煊.数据库系统概论(第5版)[M].北京:高等教育出版社.2014年
- 廖雪峰的官方网站.Python教程[OL].http://www.liaoxuefeng.com/wiki/0014316089557264a6b348958f449949df42a6d3a2e542c000.2016年
- Microsoft Virtual Academy.使用Python编程简介[OL].https://mva.microsoft.com/zh-cn/training-courses/-python--8360?l=EK9zuOO8_2604984382.2016年
- Ryan Mitchell. Python网络数据采集[M].北京:人民邮电出版社.2016年
- 宗成庆.统计自然语言处理(第2版)[M].北京:清华大学出版社.2013年
- Steven Bird, Ewan Klein, Edward Loper. Python自然语言处理[M].北京:人民邮电出版社.2014年
- Reza Zafarani, Mohammad Ali Abbasi, Huan Liu. 社会媒体挖掘[M].北京:人民邮电出版社.2015年
- 范淼,李超.Python机器学习及实践:从零开始通往Kaggle竞赛之路[M].北京:清华大学出版社.2016年
- Igor Milovanovic.Python数据可视化编程实战[M].北京:人民邮电出版社.2015年
- Toby Segaran.集体智慧编程[M].北京:电子工业出版社.2009年







The End of the Lectures

Thank You



http://www.wangting.ac.cn