

New Media Data Analytics and Application Lecture 11: 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



先定一个能达到的小日标









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 with out 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.









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 Architecture							
Weighted Av Word Cloud	verage Computing I, Media Attention	Film Box Office Prediction			Byproduct		
Historical Film Box Office Records Statistical Computing of News Report				n	News Analysis		
Flask, Word Freq	Flask, Word Frequency Computing Keyword Feature Selection Key				eyword Input		
И	/ord Dictionaries	Word Segmentation		Preprocessing			
	My SQL	Database					
上海外回语大学	Python	Web Crawlers		V	Veb APIs		





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









Text Mining

Software Development

Python PyCharm Flask MySql

Lecture	11 - [D:\SISU\新媒体数据分	析与应用\Lecture 11]	\Tag.py	- РуС	arm Community Edition 2016.2.3		
ji 👘		efactor R <u>u</u> n <u>T</u> ool	ls VC <u>S \</u>	<u>N</u> indo	/ <u>H</u> elp		
lectu	re ['lektʃə] 🕻 🚺 详细»					Tan 🛪 🕨 🕷 📿	
_ n. 演i	井; 讲稿; 教训						
🖻 vt. 滇	讲;训诫	© ≑ ¶+ It	👘 Word	SegBI	IM.py × 🛛 🎼 test.py × 🛛 🛃 Tag.py ×		
vi. 讲]祼; 讲演	☆析与应用\Lecture 1	397		Please input the Frequency of Keyword: <input name="KeyWordFrequency"/>	•	
- Forder #3	IV		398		Start Date - End Date of News <input name="StartDate"/> ånbsp;ånbsp;ånbsp;ånbsp; <ix< p=""></ix<>	nput name="EndDate">	
1732A144			399		<button type="submit">Predict</button>		
讲座	浜咞 冴咲		400		'''		
推广 -	下班刷微信不加和老外学英语		401				
1007 -			402	@app	route('/FilmBoxOffice', methods=['POST'])	_	
			403	def	filmBoxOffice():	_	
			404		: 获取当前文件路径	_	
			405		:file 为当前文件, 在ide中运行此行会报错,可改为		
			406		t d = path. dirname ('. ')		
			407		= path.dirname(file)	_	
			408				
			409		ontent = request.form['FilmWame'] #接收电影名称		
			410		tartDate = request.form['StartDate']		
			411		ndDate = request. form['EndDate']		
			412		ey#ordFrequency = request.form['KeyWordFrequency']	-	
			413				
			414		elevantFilmFrequency={}		
			415		elevantFilmBoxOffice={}		
			416		ord_tagging={}		
			417			_	
			418		· 连接到MySQL数据库		
			419		1. Connection Open	_	
			420		onn = pymysql.connect(user='root', password='123456', database='filmboxoffice', charset="utf8")		
			421		2. Cursor Creating:		
			422		ursor = conn. cursor ()	=	
			423		3. SQL Execution		
			424		·执行SQL语句,循环插入记录:		
			425		glstr = "SELECT CONTENT_WORD_SEG FROM FILM_HEWS WHERE NEWS_CONTENT LIKE '%" + content + "%' and (publ	lish_date>='"+StartDate+"' and publes	
			426		4. Cursor Moving		
			427		• 执行,游标移至当前位置		
			428		ursor. execute (sqlstr)		
Platfo	rm and Plugin Updates: Pv	Charm Community Fo	dition is re	adv to	update. (2016/12/19 3:49)	2:30 CRIF: UTE-8 % 🕀 🔿	





Home

Keyword Tagging

Keyword Comparison



Input for Keyword Comparison

Please input the Keywords: 捉妖记

西游降魔篇

Please input the Frequency of Keyword:

Start Date - End Date 2013-1-1

2016-12-1

Comparison



















← → C ③ 127.0.0.1:5000/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

← → C ③ 127.0.0.1:5000/FilmBoxOffice

<u>Home</u>

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

Film Box Office Prediction of 长城

《芳华》

Film Box Office of 芳华: 151097.2136392405(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

EXAMPLE 2: Film Stars

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:

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 《鸩之媚》

Guess for Hypothesis

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

2010-032010-122011-092012-072013-042014-012014-112015-042016-042017-032017-12

Tagging Analysis

2010-032010-122011-092012-072013-042014-072014-172015-062016-062017-032017-12

电视行业 景 米大

信葉シーに

创作

画北

现在

3.0 -

2.5 -

2.0 -

1.5 -1.0 -0.5 -

0.0

Similarity: 55.70247933884298%

0.00

2010-032010-122011-092012-022013-042014-012014-112015-062016-062017-032017-12

	明星	112
1	年代	
1	阿里	2 51
'n	自制	
- 1	视频	219 233
-	网站	
-	平台	40 306
-	内容	159 246
- 1	合作	
-	联合	3
	Blé	
	未来	
	实现	<mark>.</mark> 37
	盈利	
	版权	
	Ђ Ђ	
	想到	
	設量	
1	没有	
-	广告	
	8)/F	
	广告商	
	影视	167 525
	游戏	10 122

2010-032010-122011-092012-072013-042014-012014-112015-062016-062017-032017-12

0.0 ⊥

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: <u>http://www.public.asu.edu/~huanliu/</u>

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 Course

Thank You

http://www.wangting.ac.cn