

New Media Data Analytics and Application

Lecture 8: Quantitative Analysis for Online Journalism
Ting Wang

Outlines

- 1. Data Analysis for Online Journalism
- 2. The Foundation of Statistics
- 3. Pearson Correlation Coefficient
- 4. Bayes' Theorem
- 5. Markov Model







some analysis approaches for data journalism

Data Analysis for Online Journalism





Now, We have data.

What shall we do in the next step for online journalism analysis?



1. How to extract information from data? Data Clean and Preprocessing

2. How to measure the information of news? es Quantitative Modules

Comparison, Classification and Clustering 3. How to analysis these information?





围脖关键词

围脖关键词利用自然语言处理的关键词抽取技术,分析用户近期发表领博内容,提取代表用户兴趣的关键词,并采用文档可视化技术对关键词进行可视化,便于用户快速了解自己、好友、主题等的关键词。

使用以下账号登录





http://app.thunlp.org/



April, 2016

November, 2016



papi酱





罗辑思维退出papi酱令人深思 网红经济不行了吗?

3685

□ 我要评论

2016-11-25 15:23 来源:新京报





罗辑思维退出papi酱,网红经济不行了吗?

对个人形象的克制使用和保持健康,才是网红经济不成为一锤子买卖的关键。

据报道,罗辑思维已与著名网红papi酱分手。记者调查发现,早在今年8月29日,papi酱所在的公司春雨听雷在股东一栏里就去掉了罗辑思维的投资经营主体北京思维造物投资管理有限公司。



欢迎关注"创事记"的微信订阅号: sinachuangshiji



Technical Approaches

- 1. Keyword Extraction and Tag Analysis
- 2. News Tracking
- 3. News Alignment and Comparison
- 4. Location-based News Analysis
- 5. ...

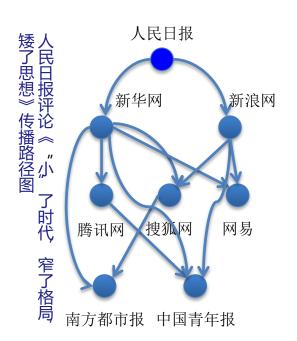


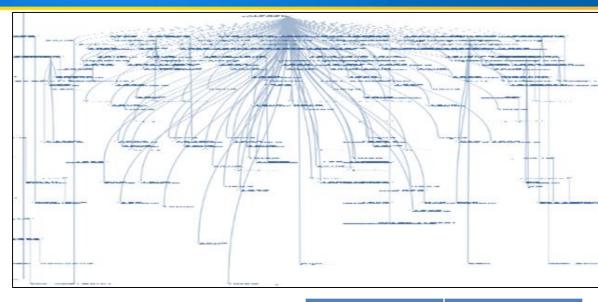
1. Keyword Extraction and Tag Analysis





2. News Tracking





Communication Efficiency Analysis base on Big Data

媒体	传播效率
新华网	71.67
新浪网	31.0
搜狐网	16.2
中国青年报	9.22
南方都市报	6.41



3. News Alignment and Comparison



《韩日军事秘密保护协定》将是韩国摆脱日本殖民统治后与日本之间签订的第一个军

事协定。如果该协议签订,韩日今后有望互通有关朝鲜军队、朝鲜社会动向、朝核以及导





China

Between **Different** Languages, Websites, Nations, and People

Notes: 1. Words in the same color have the same meaning in translation

2. The size of the word represents the importance of the word, the larger, the more important

Japan

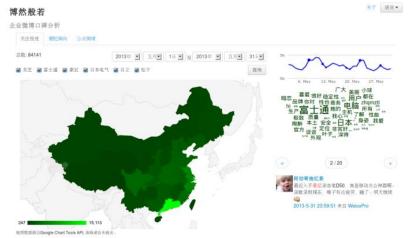
missile South Korea





4. Location-based News Analysis

- (1). Sentiment Analysis
- (2). Trend Analysis









大號: 外鱼岛等作波及日本在华大型
企业。 9月17日,包括在他、松下和
東王基在内的企业相继传来工厂生产
暂停的消息。 另指为88。 9月18
日、华富南场、优太准。 7-11等日营
零售企业也在运商大营的停业。 相关
企业表示,此类是为了保障员工的安全。 记者非晚从国际溶解和事务局
内陷于于韩晚抵法北京。 开始对中国
神经动,近日中国多城市垃圾及日示



外媒: The 81st anniversary of a Japanese invasion brought a fresh wave of anti-Japan demonstrations in China on Tuesday, with thousands of protesters venting anger over the colonial past and a current dispute involving contested islands in the East China Sea.China Warns of Further Actions as Anti-Japan Protests Resume.At least two of 11 Chinese ocean surveillance and fishery patrol ships sailing near East China Sea islets claimed by both Tokyo and Beijing have entered what Japan considers its territory.



◎ 2012 - 2013 清华大学智能技术与系统国家重点实验室信息检索组 (THUIR)



The influence by *Under the Dome*, made by Jing Chai, February 28, 2015



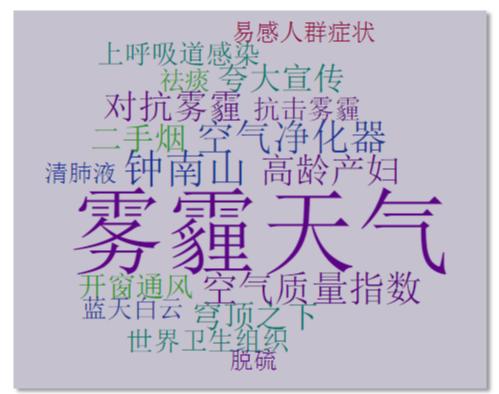
Data Description

Keyword Category	Keywords	Number of Weibo	Number of Weibo without repetition
呼吸系统疾病	呼吸系统疾病支气管炎 哮喘 咳嗽感冒胃肠型感冒、咽炎、支气管肺炎、上呼吸道感染、尘肺、结核病、鼻炎、咽喉炎、鼻窦炎、扁桃体炎	11321	6648
肿瘤	肿瘤新生儿肿瘤肺癌肺肿瘤	10655	7005
汽油	汽油质量	16	14
发电	发电、电力行业	1620	1032
净化	净化器、清新剂	9887	4770
煤	燃煤、煤炭	6587	4613
	总计	40086	24082



Keyword Extraction Based on Weibo

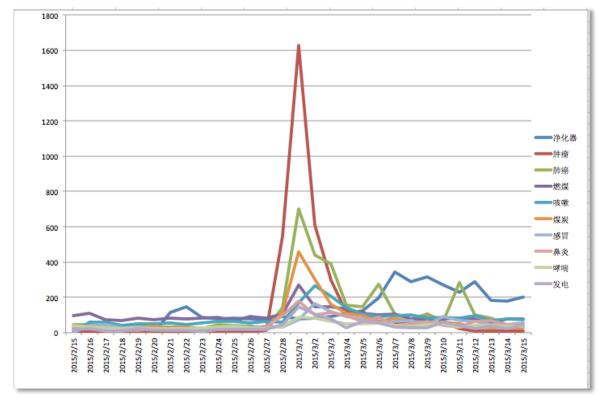
- What can you find in this graph?
- What will you do for your group?





Keyword Analysis Based on Weibo

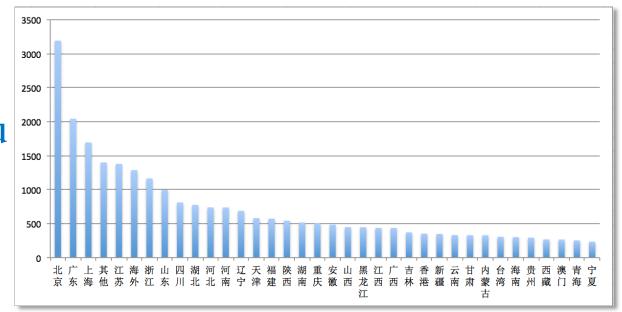
- What can you find in this graph?
- What will you do for your group?





Keyword Analysis Based on Weibo

- What can you find in this graph?
- What will you do for your group?





Conclusions:

- Keyword is an abstract of online media
- The frequency of using keywords is important



Correlative Scientific Technologies

Natural Language Processing Statistics Machine Learning ARTIFICIAL INTELLIGENCE Machine Translation Psychology Computer Science Linguistics

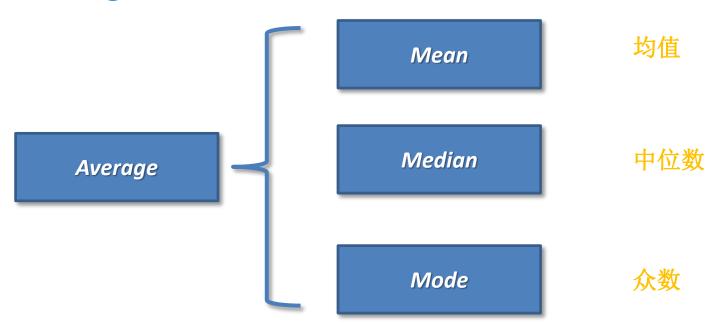




introduce some basic statistical metrics to you

The Foundation of Statistics

Average 平均数





Mean 均值

Supposing: $X=(x_1, x_2, ..., x_n)$

$$\bar{X} = \frac{\sum X}{n}$$





Median 中位数

the value separating the higher half of a data sample, a population, or a probability distribution, from the lower half.

Supposing:
$$X=(x_1, x_2, ..., x_n)$$

Sort *X* from small number to large number,

- -if *n* is an odd number, then the Median of *X* is the middle one,
- —if *n* is an even number, then the Median of *X* is the **mean** of the two middle numbers.

```
1, 3, 3, 6, 7, 8, 9

Median = 6

1, 2, 3, 4, 5, 6, 8, 9

Median = (4 + 5) ÷ 2

= 4.5
```

Mode 众数

the value that appears most often in a set of data

Comparison of common averages of values { 1, 2, 2, 3, 4, 7, 9 }

Туре	Description	Example	Result
Arithmetic mean	Sum of values of a data set divided by number of values: $ar{x}=rac{1}{n}\sum_{i=1}^n x_i$	(1+2+2+3+4+7+9) / 7	4
Median	Middle value separating the greater and lesser halves of a data set	1, 2, 2, 3 , 4, 7, 9	3
Mode	Most frequent value in a data set	1, 2, 2, 3, 4, 7, 9	2

Range 极差

the difference between the largest and smallest values

$$r = Max - Min$$





Variance 方差

the expectation of the squared deviation of a random variable from its mean, informally measures how far a set of (random) numbers are spread out from their mean, also known as D(X), Var(X)

$$s^2 = \frac{\sum (X - \bar{X})^2}{n - 1}$$

Why n-1?





Standard Deviation 标准差

$$s = \sqrt{\frac{\sum (X - \bar{X})^2}{n - 1}}$$







Expected Value 数学期望

$$E[X] = \bar{X} = \sum_{i} x_i P_i$$

Where: P_i is the weight of x_i in Statistics, P is the probability.



Properties of Expected Value

- If C is a constant, E[C]=C
- If *X* and *Y* are random variables such that $X \le Y$, then $E[X] \le E[Y]$
- -E[X+C]=E[X]+C
- -E[X+Y]=E[X]+E[Y]
- -E[CX]=CE[X]
- $-D[X]=E[X^2]-(E[X])^2$





Covariance 协方差

a measure of the joint variability of two random variables

$$Cov(X,Y) = E[(X - E[X])(Y - E[Y])]$$

= $E[XY] - 2E[Y]E[X] + E[X]E[Y]$
= $E[XY] - E[X]E[Y]$

$$S^{2} = \frac{\sum (X - X)^{2}}{n - 1}$$



Properties of Covariance

$$D(X+Y)=D(X)+D(Y)+2Cov(X, Y)$$

$$D(X-Y)=D(X)+D(Y)-2Cov(X, Y)$$

$$Cov(X, Y)=E(XY)-E(X)E(Y)$$

$$Cov(X, Y)=Cov(Y, X)$$

 $Cov(aX, bY)=abCov(X, Y)$
 $Cov(X_1+X_2, Y)=Cov(X_1, Y)+Cov(X_2, Y)$





Uncorrelatedness and independence

• If *X* and *Y* are independent, then their covariance is 0.

$$E[XY]=E[X]E[Y]$$

• The converse, however, is not generally true.

When the covariance is *normalized*, one obtains the Pearson correlation coefficient, which gives the goodness of the fit for the best possible linear function describing the relation between the variables. In this sense covariance is *a linear* gauge of dependence.

https://www.zhihu.com/question/20852004





a measurement of correlation

Pearson Correlation Coefficient

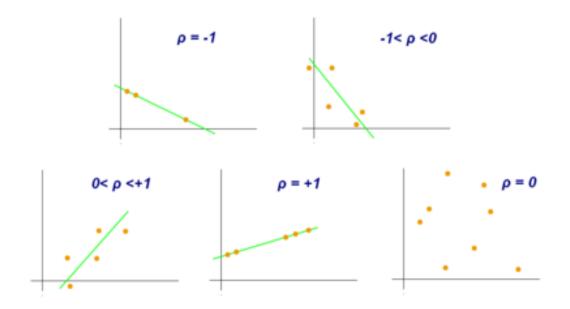
Pearson Correlation Coefficient

Pearson Correlation Coefficient is a measure of the linear correlation between two variables *X* and *Y*.

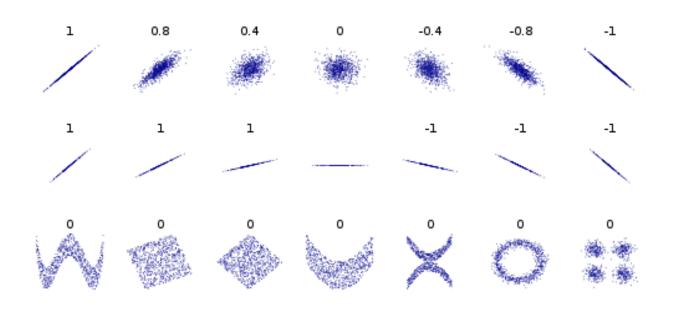
$$r(X,Y) = \frac{Cov(X,Y)}{\sqrt{Var[X] Var[Y]}}$$

- It has a value between +1 and -1, where 1 is total positive linear correlation, 0 is no linear correlation, and -1 is total negative linear correlation.
- It is widely used in the sciences. It was developed by Karl Pearson from a related idea introduced by Francis Galton in the 1880s.

Examples of scatter diagrams with different values of correlation coefficient (ρ)









Calculate PCC in Python

```
import numpy as np
```

a=np. array([1, 2, 3, 4]) b=np. array([8, 7, 6, 5])

print(np. corrcoef(a, b))





EXAMPLE 3: INOHERB







EXAMPLE 4: GARLSBERG

Probably the best beer in the world

For more visit carlsberg.com



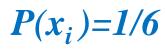


very useful for natural language processing

Bayes' Theorem

Probability 概率





Sample Space:

$$\{1, 2, 3, 4, 5, 6\}$$



$$P(x_i) = 1/2$$

$$\{H, T\}$$



Properties of Probability

$$P(x_i) \geq 0$$

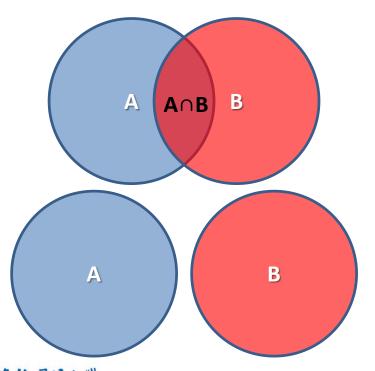
$$P(x_i) \in [0,1]$$

$$\sum_{i=1}^{n} P(x_i) = 1$$





Independence 独立性



Dependent

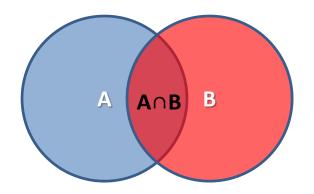
Independent



Conditional Probability 条件概率

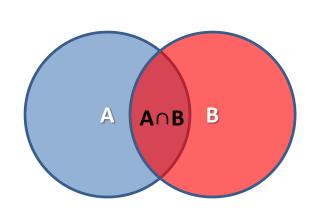
P(A | B), is the probability of observing event A given that B is true

$$P(A|B) = P(A \cap B)/P(B)$$





Bayes' Theorem 贝叶斯定理



$$P(A|B) = P(A \cap B)/P(B)$$

$$P(A \cap B) = P(A|B)P(B)$$

$$P(A \cap B) = P(B|A)P(A)$$

$$P(A|B)P(B) = P(B|A)P(A)$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$



Bayes' Theorem plays an very important role in statistical NLP.

- We can predict what you will say!
 - Uncle Sam: How are you?
 - Chinese student: Fine, Thank you, and you?
 - Chinese student's Predictive Answer: I am fine, too!
 - Uncle Sam: Nothing much.
 - Chinese student:。。。(不多??)







• Because, for Chinese students:

```
P(Fine, Thank you, and you? | How are you?)
P(I am fine, too! | Fine, Thank you, and you?)
P(Nothing much | Fine, Thank you, and you?)
```

In the corpus of Chinese students,

P(I am fine, too! | Fine, Thank you, and you?)>P(Nothing much | Fine, Thank you, and you?)



Another Example:

I ate a red _____

A. telephone B. light C. swim D. tomato



No Grammar! But the Frequency of use!

• The most successful Chinglish: Long time no see!

• Chinglish Future Star: Good Good Study, Day Day UP!





your future is decided by now, not the past

Markov Model

Stochastic Process 随机过程 Markov Chain 马尔科夫链



$$X=(x_1, x_2, \ldots, x_n)$$

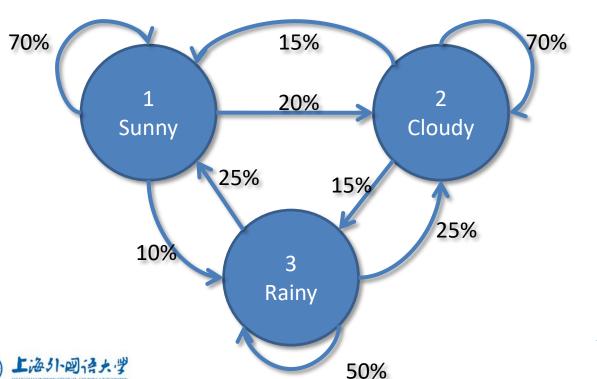
x, is a Stochastic Process

1,3,5,2,1,4,2,6,3,.....

X is a Markov Chain



Transition Probability 转移概率



$$\begin{bmatrix} P_{11} & P_{12} & P_{13} \\ P_{21} & P_{22} & P_{23} \\ P_{31} & P_{32} & P_{33} \end{bmatrix} =$$

$$egin{array}{cccc} 0.7 & 0.2 & 0.1 \ 0.15 & 0.7 & 0.15 \ 0.25 & 0.25 & 0.5 \ \end{array}$$

Stochastic Matrix 概率转移矩阵

出度之和100%

Markov Model 马尔科夫模型

$$P(x_{t+1}|x_1, x_2, \dots, x_t) = P(x_{t+1}|x_t)$$

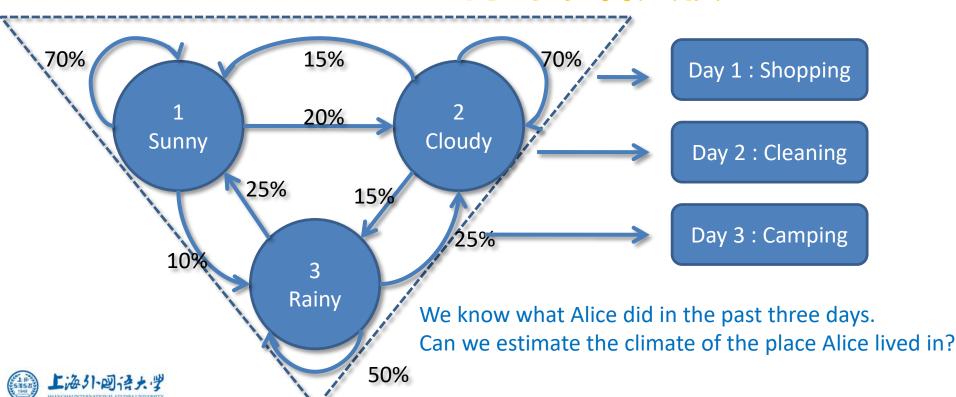
First-Order Markov Model

Your future is not decided by your past, but now!

Second-Order Markov Model

$$P(x_{t+1}|x_1, x_2, \dots, x_t) = P(x_{t+1}|x_t x_{t-1})$$

Hidden Markov Model 隐马尔科夫模型



The Applications of Markov Model in NLP

- Machine Translation
- Word Segmentation
- Speech Recognition
- Part-of-speech Tagging
- Natural Language Generation









Reference

Reference

• https://item.jd.com/11701113.html

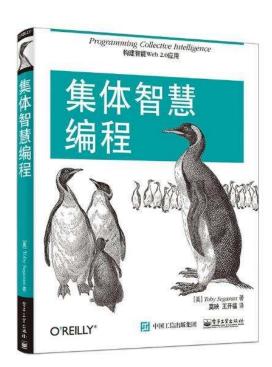






Reference

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









The End of Lecture 8

Thank You

http://www.wangting.ac.cn

