

### Computational Journalism

Lecture 1: A Brief Introduction

Ting Wang

### Outlines

- 1. Significance of Computational Journalism
- 2. Definition of Computational Journalism
- 3. History of Computational Journalism
- 4. Domains of Computational Journalism





the significance of computational journalism

### Why Data Analysis





US Dollars,
The King of
the World





Bretton Woods, 1944



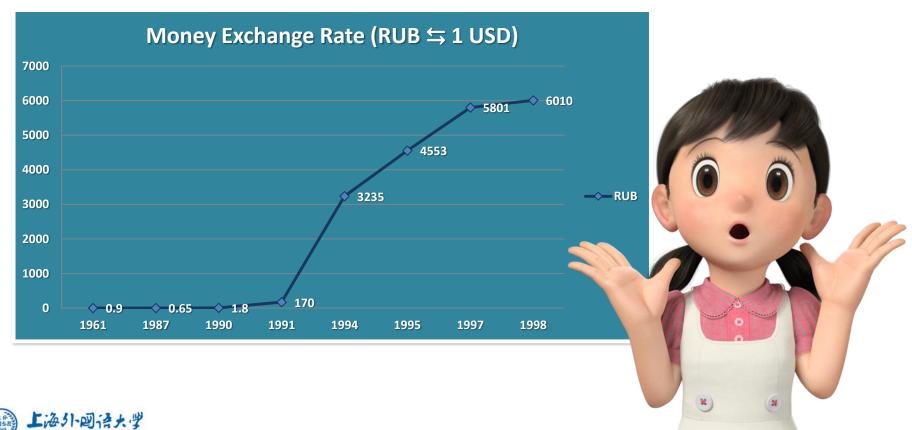
Political events impact on economic trends













### http://cn.investing.com/







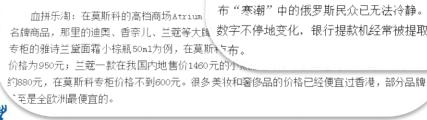


小柒海外正品代购 🥏 🙏

扫一扫上面的二维码图案,加我强信

#### **卢布贬值引发代购潮,俄罗斯华人大淘便宜货**

由于沂期卢布的连续贬值,用 消售范围主要围绕化妆、日月







#### News and New Media are crucial to the investment.









网易首页 应用 ~

网易考拉 ~

### 英国脱欧民意暂领先3%(组图)



资讯

凤凰网资讯 >滚动新闻 >正文

距离英国公投还剩4天 "留欧"民意支持率反超"脱欧"

2016年06月20日 04:41

来源:新快报

鳳凰 资讯

凤凰网资讯 >滚动新闻 >正文

美媒:英国"脱欧"公投很可能失败 民调靠不住

2016年06月22日 00:11

来源:参考消息网



















# Real time data analysis based on news and my EURO currency exchange

http://finance.qq.com/zt2016/gongtou/index.htm

1	
News	Rate
Gibraltar (IN)	1:7.52
New Castle (IN)	1:7.51
Sunderland (OUT)	1:7.47
Oxford (IN), North Ireland (IN)	1:7.43
49.79% (IN), 50.21% (OUT)	1:7.30
48.95% (IN), 51.05% (OUT)	1:7.25
Wales (OUT)	1:7.22
48.28% (IN), 51.72% (OUT), 339/382 Regions	1:7.24
	News Gibraltar (IN) New Castle (IN) Sunderland (OUT) Oxford (IN), North Ireland (IN) 49.79% (IN), 50.21% (OUT) 48.95% (IN), 51.05% (OUT) Wales (OUT)





Final Currency Exchange Results:

(7.52-7.24)\*2200=616 RMB







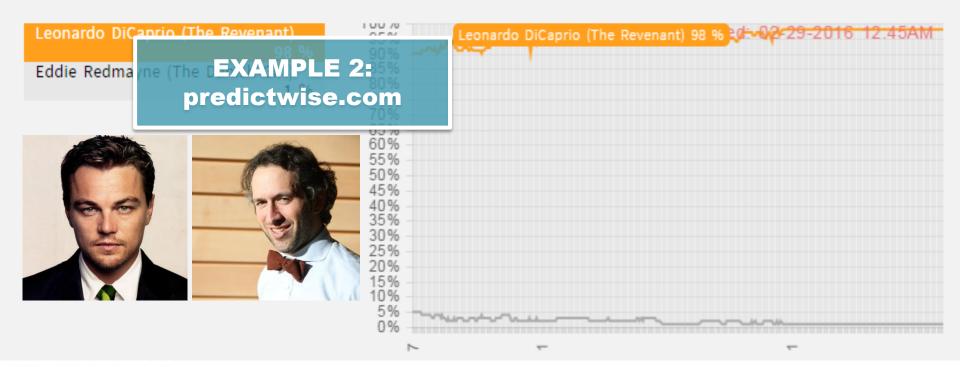




Do you have a similar experience using data analysis based on news?



#### **Leading Actor**







David Rothschild,

PhD of Wharton School of Business at the University of Pennsylvania Microsoft researcher at Microsoft Research in New York City He correctly predicted 50 of 51 Electoral College outcomes in February of 2012, average of 20 of 24 Oscars from 2013-5, and 15 of 15 knockout games in the 2014 World Cup.

- POLITICS
- SPORTS
- ENTERTAINMENT
- ECONOMIC/FINANCIAL



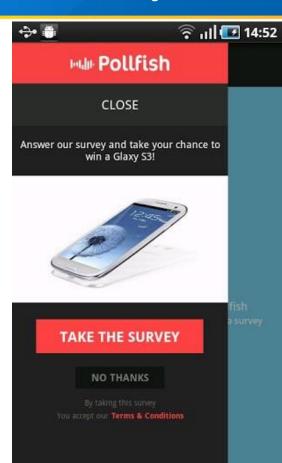
### **Approaches**

Data Collection:Pollfish, MSN, Xbox

 Data Analysis:
 Statistical Analysis based on Historical Data

http://predictwise.com/



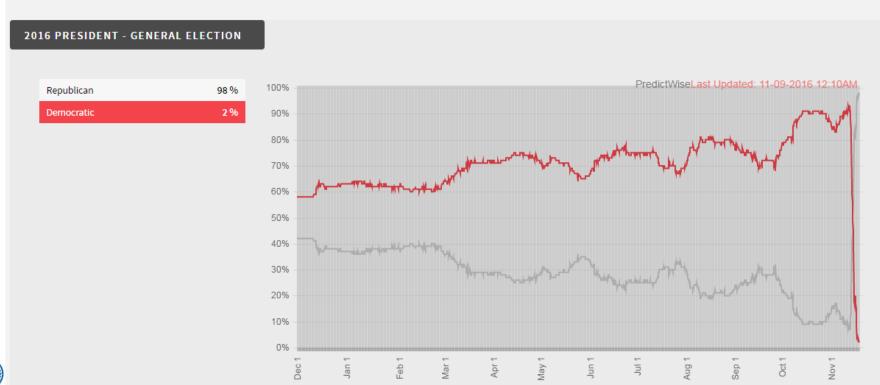


#### **Politics**

#### **2016 PRESIDENT - GENERAL ELECTION** PredictWise Last Updated: 09-05-2016 1:40PM 90% 75 % Democratic 85% 25 % Republican 80% 75% 70% 65% 60% 55% 50% 45% 40% 35% 30% 25% 15% 10% Aug 1



#### **Politics**





1. Examine data quality - in this election polls were not reaching all likely voters

2. Beware of your own biases: many pollsters were likely Clinton supporters and did not want to question the results that favored their candidate. For example, Huffington Post had forecast 98% chance of Clinton Victory.





# What is your opinion?

https://markets.predictwise.com/politics/trump-specials?q=/politics/&

**Predict**Wise

Market data here! But, PredictWise stoked to expand from collecting & analyzing market of survey data (and beyond) for public opinion. Evergreen, detailed time-trends on public opinion, and how stakeholders engage with people earned media.

ABOUT BLOG MARKETS RESEARCH

#### **Trump Specials**



When will Trump stop being President of USA?

Outcome	Market	Derived Betfair Price	Betfair Back	Betfair Lay	Derived Predictlt Price
2020 to 1/20/2021	5 %	\$ 0.744	1.34	1.35	
Finishes 1st Term	70 %	\$ 0.697	1.43	1.44	
2019	20 %	\$ 0.200	4.90	5.10	\$0.230
2018	5 %	\$ 0.053	18.00	19.50	\$ -0.440
2017	0 %	\$ 0.001	470.00	0.00	\$ 0.505



Last Updated: 10-07-2018 12:17PM



the definition of computational journalism

### What is Computational Journalism

### The significance of Data Analysis

- 1. To obtain new information
- 2. To enlarge the benefits
- 3. To avoid the risks





# INFORMATION DISCOVERY CONCLUSION SUGGESTING PECISION SUPPORT

are three objectives of data analysis



### Definition:

https://en.wikipedia.org/wiki/Computational\_journalism

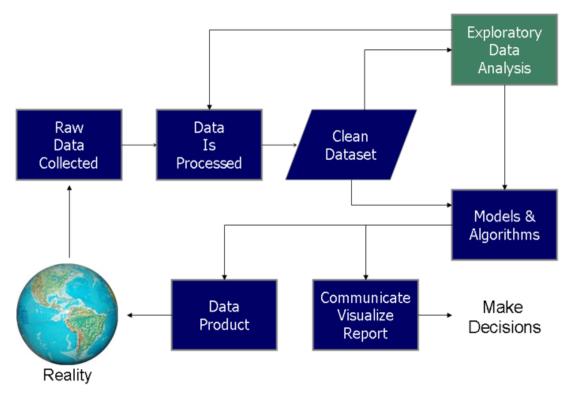
Computational Journalism can be defined as the application of computation to the activities of journalism such as information gathering, organization, sensemaking, communication and dissemination of news information, while upholding values of journalism such as accuracy and verifiability.



- Two types of decisions:
  - Quantitative Decision with a value
    - Prediction, Regression

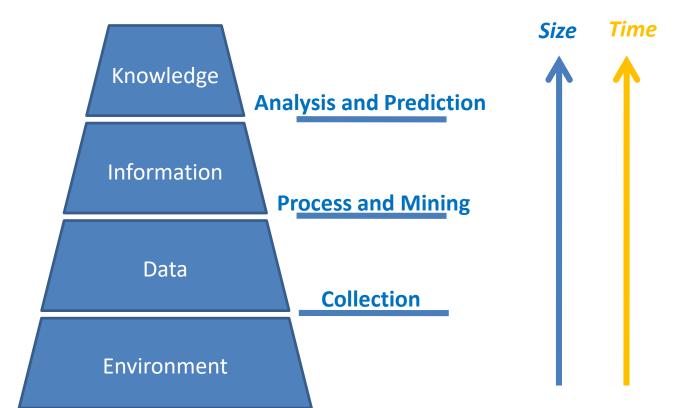
- Qualitative Decision with a label
  - Classification, Clustering

#### **Data Science Process**





Relationship between data, information and knowledge





Methodology

Data Cleaning and

**Information Retrieval** 

**Structured** 

Web

**Information** 

**Knowledge Fusion and Information Updating** 

> **Feature Analysis** and **Updating**

**Personalized Products or** 

**Services** 

**Prediction or Classification** 



Web Data

**Information** 

**Acquisition** 















Do you have some methods to forecast the auto sale market in the next several months?



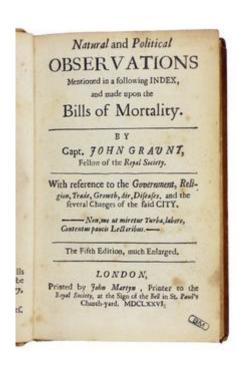




computational journalism will be everywhere in the future

### The History of Computer Data Analytics

### The History of Computer Data Analytics



# 社会统计1662

John Graunt (24 April 1620 – 18 April 1674) used statistical analysis to predict the onset and spread of bubonic plague in London, which led him to the Royal Society.



### The History of Computer Data Analytics



Thomas Bayes (1701-1761)

## 贝叶斯决策理论1763

Bayes, Thomas; Price, Mr. (1763). "An Essay towards solving a Problem in the Doctrine of Chances. 《机会问题的解法》". Philosophical Transactions of the Royal Society of London. 53 (0): 370–418. doi:10.1098/rstl.1763.0053



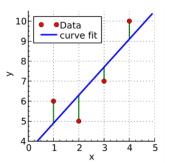
### The History of Computer Data Analytics

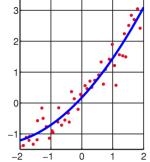


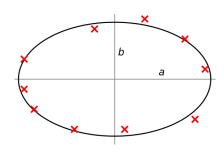
Carl Friedrich Gauss (1777–1855)

# 最小二乘法1805

Least squares for data fitting and regression







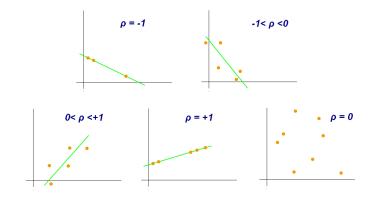




Karl Pearson (27 March 1857 – 27 April 1936)

# 相关系数1880

#### Pearson Correlation Coefficient



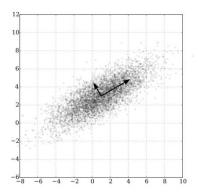




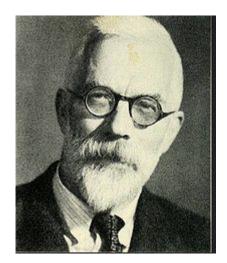
Karl Pearson (27 March 1857 – 27 April 1936)

# 主成分分析1901

#### Principal Component Analysis



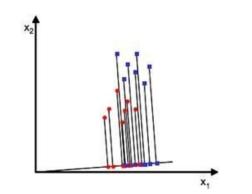


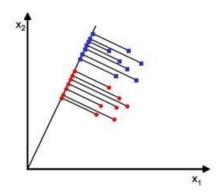


R. A. Fisher (17 February 1890 – 29 July 1962)

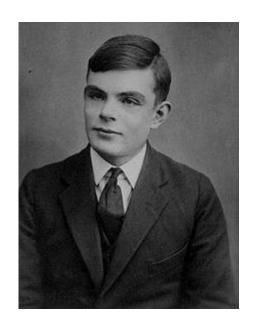
# 线性判别 1936

#### Linear Discriminant Aanalysis





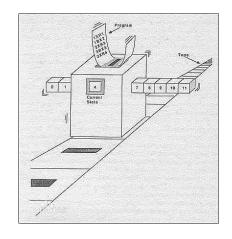




Alan Turing (23 June 1912 – 7 June 1954)



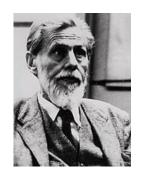
**Turing Machine** 







Warren McCulloch

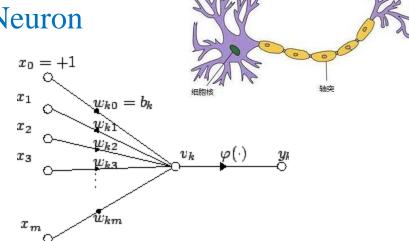


Walter Pitts

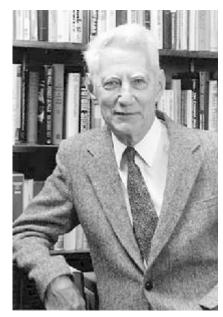
# 人工神经元1943

#### **Artificial Neuron**

$$y_k = arphi \left( \sum_{j=0}^m w_{kj} x_j 
ight)$$







Claude Shannon (April 30, 1916 – February 24, 2001)

# 信息论1948

Information theory

- -Entropy (信息熵)
- -Mutual Information (互信息)



2006, 50 Years Anniversary



# 人工智能1956

Summer Research Project on Artificial Intelligence, Dartmouth College





Prof. Pedro Domingos University of Washington

#### 2015, ACM

- 1. 符号主义
- 2. 联结主义
- 3. 进化主义
- 4. 贝叶斯主义
- 5. 类推主义

#### Five Tribes in AI

Tribe	Origins	Master Algorithm
Symbolists	Logic, philosophy	Inverse deduction
Connectionists	Neuroscience	Backpropagation
Evolutionaries	Evolutionary biology	Genetic programming
Bayesians	Statistics	Probabilistic inference
Analogizers	Psychology	Kernel machines





Avram Noam Chomsky



**Herbert Simon** 



Allen Newell

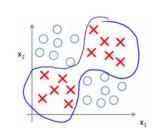
# 符号主义 1957 Symbolism

- 1. Plato is a man.
- 2. Man will die.
- 3. Plato will die.

Expert System
Universal Grammar and Chomsky Hierarchy







1950

1960

# 联结主义 1943 connectionism

2000

8层袖经网络

2020

2010

**Donald Olding Hebb** 







1980

1990

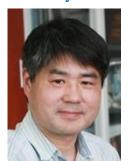
1970

Frank Rosenblatt





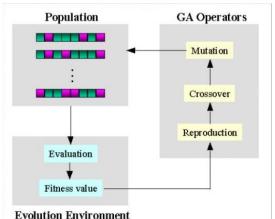
John Henry Holland



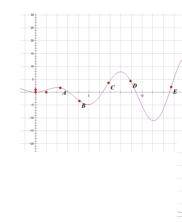
Yuhui Shi

# 選化主义 1970's Evolutionism

Genetic Algorithm



Particle Swarm Optimization

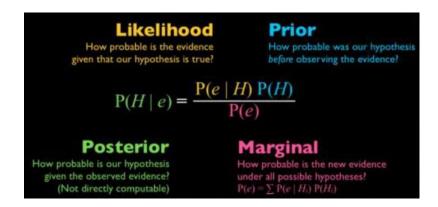






Judea Pearl

# 贝叶斯主义 1763 Bayesianism





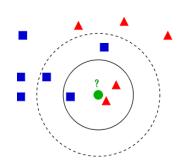


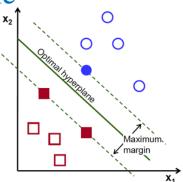
Vladimir Vapnik

# 类推主义 1951 Analogism

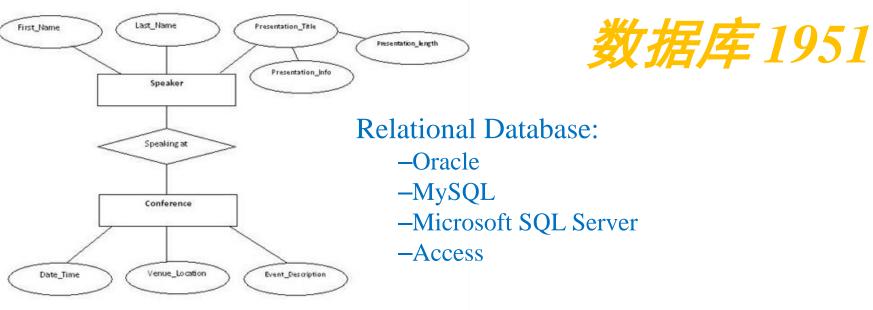
K-Nearest Neighbour

Support Vector Machine







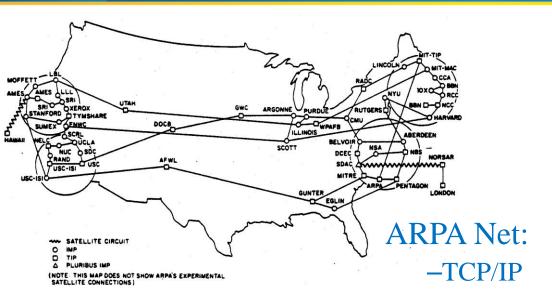


**Entity Relationship Diagram** 

NoSQL:

- MongoDb





互联网1969

#### The Internet:

- -World Wide Web
- -Email







Web2.0 2004

Web 1.0 Contents made by Providers Web 2.0 Contents made by Customers

Semantic Web: Web 3.0?





# 物联网2005

#### The Internet of Things





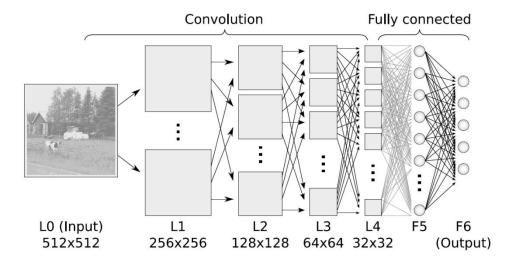




Geoffery Hinton

# 深度学习 2006

Deep Learning







# 云计算2006

#### **Cloud Computing**

- -Grid Computing
- -Distributed Computing
- -Parallel Computing
- -Utility Computing

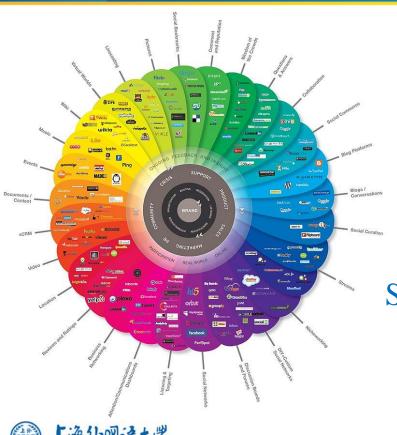
**—**…











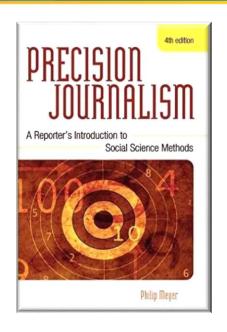
# 社交媒体的繁荣

2000's

Social Media







#### 数据新闻1967

#### Data Analysis in Journalism

Philip Meye, Precision Journalism, 1970's

The National Institute for Computer-Assisted Reporting – NICAR, 1994

Computational Journalism, Georgia Institute of Technology, 2006

Frontiers of Computational Journalism, Columbia Journalism School, 2012

Masters in Computational Journalism, Syracuse University, 2015

Computational Journalism Lab, Stanford University, 2015





knowledge domains of computational journalism

## Domains of New Media Data Analytics

#### Domains of New Media Data Analytics

#### Relevant Disciplines

- Journalism
- Computer Science
- Mathematics
- Psychology
- Economics
- Politics
- Linguistics
- **—** ....





#### Domains of New Media Data Analytics

#### Corresponding Technologies

- Computer Science
- Artificial Intelligence
- Machine Learning
- Statistical Analysis
- Natural Language Processing
- Pattern Recognition
- **—** ...



#### Domains of New Media Data Analytics

#### Related Fields

- 1. Database journalism
- 2. Computer-assisted reporting
- 3. Data-driven journalism





#### Home Work

#### Home Work

1. Identify at least three major side effects of information sharing on social media.

2. Rumors spread rapidly on social media. Can you think of some methods to block the spread of rumors on social media?







#### The End of Lecture 1

Thank You

http://www.wangting.ac.cn

