# New Media
# Data Analytics and Application

## Lecture 1: A Brief Introduction

Ting Wang

1. Significance of Data Analysis

2. Definition of Data Analysis

3. History of Computer Data Analytics

4. Domains of New Media Data Analytics

the significance of data analysis

# Why Data Analysis

EXAMPLE 1:
Money Exchange

US Dollars,

The King of the World



*Bretton Woods, 1944*

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

## Political events impact on economic trends

Money Exchange Rate (RUB ⇆ 1 USD)

卢布贬值引发代购潮，俄罗斯华人大淘便宜货

2014-12-10 10:10:10

个人信息

MTS RUS 12:24 100%

CHANEL

小柒海外正品代购

俄罗斯

扫一扫上面的二维码图案，加我微信

腾讯教育　高考　考研　公务员　校园　留学　外语　中小学

卢布贬值震惊全球 中国留学生做起奢侈品

留学　浙江在线-钱江晚报[微博] 2014-12-18 09:35　我要分享▾

浙江在线12月18日讯：（钱江晚报记者 钟卉）在2014年的最后
场空前巨大的危机中。卢布在本周一暴跌超10%，创下1998年来最大
斯央行一口气加息650个基点，从10.5%提高至17%，仍拦不住卢布的
贬值25%。

受国际油价持续走低、西方对俄施加经济制裁、俄国际储备处
合作用，今年以来卢布贬值幅度已经超过50%。危机当前，一惯强
布"寒潮"中的俄罗斯民众已无法冷静。俄罗斯街头出现这样的坏
数字不停地变化，银行提款机经常被提取一空，俄罗斯人冲进商场

网易新闻　网易首页 ▸ 新闻中心 ▸ 滚动新闻

卢布人贬值 快去俄罗斯海

2014-12-28 02:35:41　来源：华商报(西安)

又到岁末，人民币依然维持全年坚挺走势
是著名的金砖国家俄罗斯的卢布跌幅更是大。
少呢。

血拼乐淘：在莫斯科的高档商场Atrium
名牌商品，那里的迪奥、香奈儿、兰蔻等大牌
专柜的雅诗兰黛面霜小棕瓶50ml为例，在莫斯科卢布
价格为950元；兰蔻一款在我国内地售价1460元的小棕
约880元，在莫斯科专柜价格不到600元。很多美妆和奢侈品的价格已经便宜过香港，部分品牌
至是全欧洲最便宜的。

由于近期卢布的连续贬值，用
销售范围主要围绕化妆、日用

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

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

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

# Why Data Analysis

June 23, 2016

FUSION 2016
Heidelberg  July 5-8

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

# Why Data Analysis

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

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

| Time | News | Rate |
|------|------|------|
| 06-24 06:17 | Gibraltar (IN) | 1:7.52 |
| 06-24 06:31 | New Castle (IN) | 1:7.51 |
| 06-24 07:28 | Sunderland (OUT) | 1:7.47 |
| 06-24 09:23 | Oxford (IN), North Ireland (IN) | 1:7.43 |
| 06-24 10:01 | 49.79% (IN), 50.21% (OUT) | 1:7.30 |
| 06-24 11:07 | 48.95% (IN), 51.05% (OUT) | 1:7.25 |
| 06-24 11:40 | Wales (OUT) | 1:7.22 |
| 06-24 12:13 | 48.28% (IN), 51.72% (OUT), 339/382 Regions | 1:7.24 |



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

*Final Currency Exchange Results:*

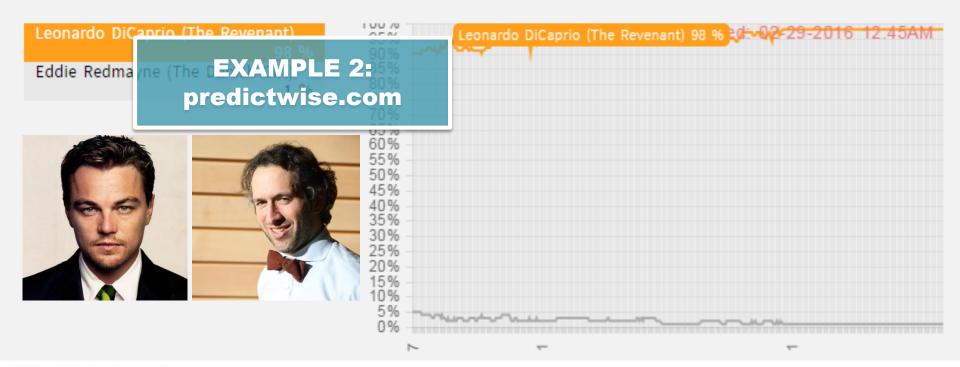(7.52-7.24)*2200=616 RMB

**Ask A Question**

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

# Why Data Analysis

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
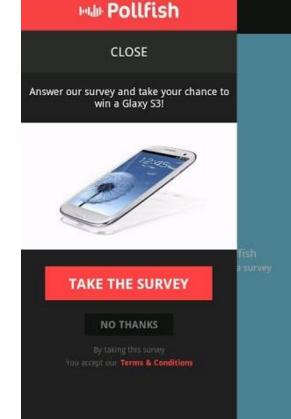- SPORTS
- ENTERTAINMENT
- ECONOMIC/FINANCIAL

## *Approaches*

– Data Collection:
  Pollfish, MSN, Xbox

– Data Analysis:
  Statistical Analysis based on
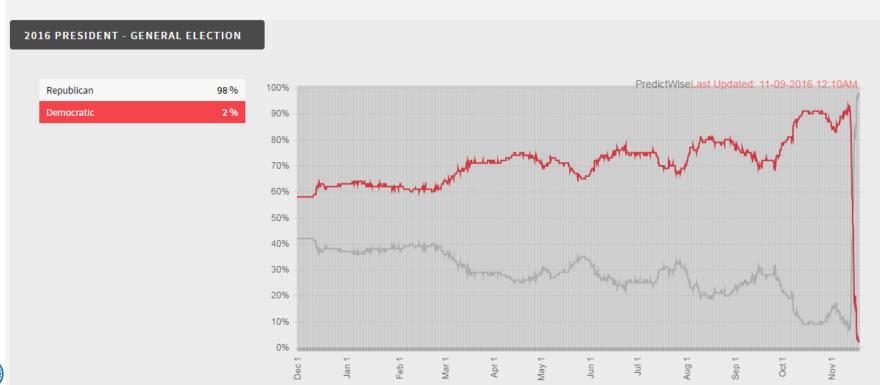  Historical Data

http://predictwise.com/



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

## Politics

**2016 PRESIDENT - GENERAL ELECTION**

| | |
|---|---|
| Democratic | 75 % |
| Republican | 25 % |



PredictWise Last Updated: 09-05-2016 1:40PM

Democratic 63% (01-11-2016 12:10PM)

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY
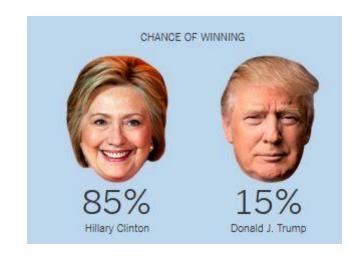
# Why Data Analysis

## 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.

CHANCE OF WINNING

85%
Hillary Clinton

15%
Donald J. Trump

http://www.kdnuggets.com/2016/11/trump-shows-limits-prediction.html

the definition of data analysis for journalism

# What is Data Analysis

*The significance of Data Analysis*

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

# INFORMATION DISCOVERY
# CONCLUSION SUGGESTING
# DECISION SUPPORT

## are three objectives of data analysis

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

**Definition:**

https://en.wikipedia.org/wiki/Data_analysis

***Analysis of data** is a process of inspecting, cleaning, transforming, and modeling data with the goal of discovering useful information, suggesting conclusions, and supporting decision-making.*
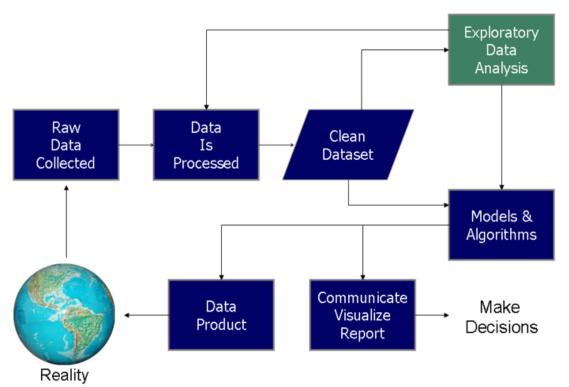
上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

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

    - Qualitative Decision with a label
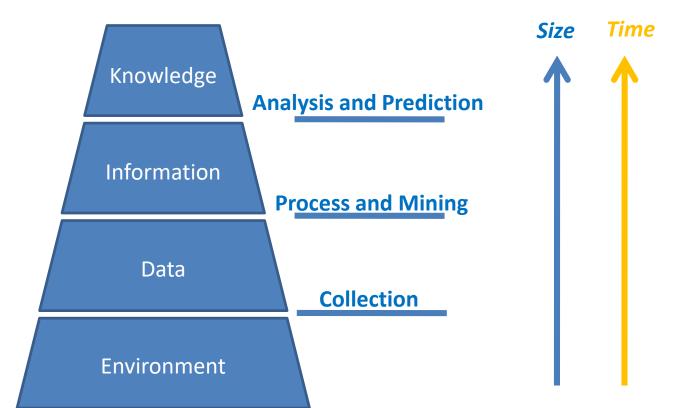        - Classification, Clustering

Data Science Process

# What is Data Analysis

*Relationship between data, information and knowledge*

## *Methodology*

**1**   **2**   **3**   **4**

*Information Acquisition*

*Data Cleaning and Information Retrieval*

*Knowledge Fusion and Information Updating*

*Prediction or Classification*

**Web Data** → **Structured Web Information** → **Feature Analysis and Updating** → **Personalized Products or Services**

围脖关键词

围脖关键词利用自然语言处理的关键词技术技术和关键词云技术技术对用户历史发表微博内容，提取代表用户兴趣的关键词，并采用可视化技术可视化，便于用户快速了解自己、好友、主题等的关键词。

**EXAMPLE 3:**
**papi酱**

http://app.thunlp.org/

papi酱美女

# What is Data Analysis



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

papi酱美女

## *Weibo Key Words for papi酱*

**1**

**2**

**3**

**4**

**Information Acquisition**

**Data Cleaning and Information Retrieval**

**Knowledge Fusion and Information Updating**

**Prediction or Classification**

*Weibo Data from Sina API*

*Word Segmentation, Part-of-speech Tagging, Translation-based and Frequency-based Key Word Extraction,…*

*Key Word List and Updating*

*Video Recommendation, Personal Current Status Analysis*

**Ask A Question**

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

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

computational journalism will be everywhere in the future

# The History of Computer Data Analytics

Natural and Political
**OBSERVATIONS**
Mentioned in a following INDEX,
and made upon the
**Bills of Mortality.**

BY
Capt. *JOHN GRAUNT*,
Fellow of the *Royal Society*.

With reference to the *Government*, *Reli-gion*, *Trade*, *Growth*, *Air*, *Diseases*, and the several Changes of the said CITY.

*Non me ut miretur Turba, laboro, Contentus paucis Lectoribus.*

The Fifth Edition, much Enlarged.

LONDON,
Printed by *John Martyn*, Printer to the *Royal Society*, at the Sign of the Bell in St. *Paul's* Church-yard. MDCLXXVI.

**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.

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
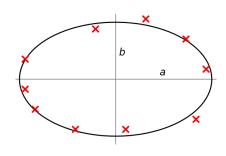
最小二乘法 *1805*

Carl Friedrich Gauss (1777–1855)
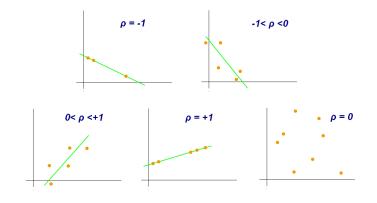
Least squares for data fitting and regression

Karl Pearson
(27 March 1857 – 27 April 1936)
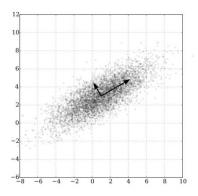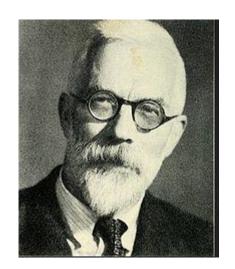
相关系数 *1880*

Pearson Correlation Coefficient



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

Karl Pearson
(27 March 1857 – 27 April 1936)

主成分分析*1901*

Principal Component Analysis



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

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

**线性判别** *1936*

Linear Discriminant Aanalysis

图灵机 *1936*

Turing Machine



Alan Turing
(23 June 1912 – 7 June 1954)

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

**Warren McCulloch**

**Walter Pitts**

人工神经元 *1943*

Artificial Neuron

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

树突

轴突末梢

细胞核

轴突

$x_0 = +1$

$x_1$

$x_2$

$x_3$

$x_m$

$w_{k0} = b_k$

$w_{k1}$

$w_{k2}$

$w_{k3}$

$w_{km}$

$v_k$

$\varphi(\cdot)$

$y_k$

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

**信息论 1948**

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

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

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

2006, 50 Years Anniversary



**人工智能*1956***

Summer Research Project on Artificial Intelligence, Dartmouth College

*Five Tribes in AI*

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

| 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 |

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY
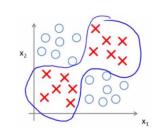
Avram Noam Chomsky



Herbert Simon



Allen Newell

*符号主义 1957*
*Symbolism*

1. Plato is a man.
2. Man will die.
3. Plato will die.
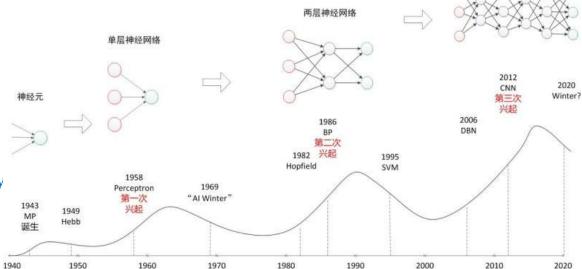
Expert System
Universal Grammar and Chomsky Hierarchy

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

Donald Olding Hebb

Frank Rosenblatt    Marvin Lee Minsky

联结主义 *1943*
*connectionism*

多层神经网络

两层神经网络

单层神经网络

神经元

2012
CNN
第三次
兴起

2020
Winter?

2006
DBN

1986
BP
第二次
兴起

1982
Hopfield

1995
SVM

1958
Perceptron
第一次
兴起

1969
"AI Winter"

1943
MP
诞生

1949
Hebb

1940    1950    1960    1970    1980    1990    2000    2010    2020

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

# The History of Computer Data Analytics

进化主义 1970's
Evolutionism

John Henry Holland

Yuhui Shi

Genetic Algorithm

Particle Swarm Optimization

上海外国语大学
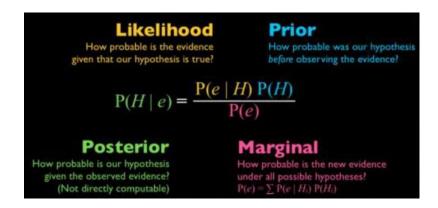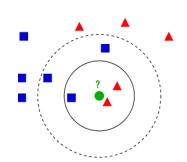SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

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

**数据库 1951**

Relational Database:
  – Oracle
  – MySQL
  – Microsoft SQL Server
  – Access

NoSQL:
  – MongoDb

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

# The History of Computer Data Analytics



**互联网 1969**

ARPA Net:
–TCP/IP

The Internet:
–World Wide Web
–Email



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

*Web2.0 2004*

Web 1.0 Contents made by Providers

Web 2.0 Contents made by Customers

Semantic Web: Web 3.0?

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

物联网 *2005*

The Internet of Things



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

深度学习 *2006*

Deep Learning



Geoffery Hinton

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

**云计算** *2006*

Cloud Computing
- Grid Computing
- Distributed Computing
- Parallel Computing
- Utility Computing
- …

大数据 *2008*

社交媒体的繁荣

*2000's*

Social Media



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY
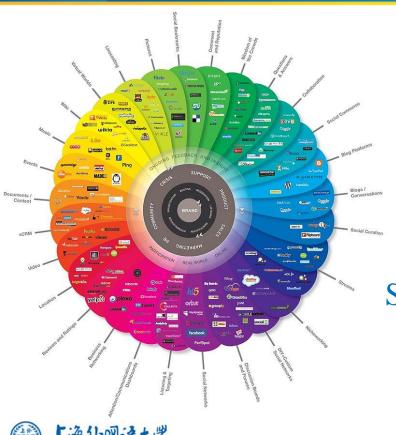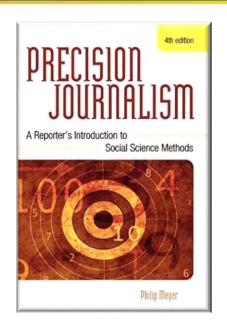
数据新闻1967

## 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

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

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
- ….

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

# Domains of New Media Data Analytics

***Corresponding Technologies***

- Computer Science
- Artificial Intelligence
- Machine Learning
- Statistical Analysis
- Natural Language Processing
- Pattern Recognition
- …

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

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?

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

# The End of Lecture 1

Thank You

http://www.wangting.ac.cn