# Making Big Profits from Big Data
**Opportunities and Challenges in Business**

Personalized ads based on **YOUR** data

Enhanced click-through rate = MORE profits!

境外消费TOP 10

1.香港  6.新加坡
2.俄罗斯  7.加拿大
3.美国  8.澳门
4.中国台北  9.巴西
5.澳大利亚  10.马来西亚

(*地区排名与地图非对应关系)

海淘热门商品TOP 5

手机　奶粉　箱包　女裙　面膜

(*海淘数据来自好奇心日报)

2014年双十一交易额

| 1分11秒 | 14分02秒 | 38分28秒 | 7时左右 | 7时17分 | 13时31分 | 15时29分 | 21时12分 | 23时59分 |
|---|---|---|---|---|---|---|---|---|
| 1亿 | 50亿 | 100亿 | 191亿 超2012年全年交易额 | 200亿 | 362亿 超2013年全年交易额 | 400亿 | 500亿 | 571.1亿 |
| 1亿 | | | | 100亿 | | | | 243.3亿 |

支付交易比例

2013: 75.97% / 24.03%
2014: 57.4% / 42.6%

PC端占比
无线端占比

无线交易增长

113 亿元 (2013) → 243.3 亿元 (2014)

微信号:imciow

- **Serving 10 million+ users every day, 36 million+ users at the busiest day**

- **40+ recommended commodities for each user**

- **Predicting user preferences from behaviors data**

**Due to better models, the revenue of Tmall.com increased by 20%+**

**User behavior data from Tmall.com (Alibaba)**

| User | Brand | Date | Behavior |
|------|-------|------|----------|
| Alice | Lenovo | 2014-04-18 | Click |
| Bob | Sony | 2014-04-20 | Click |
| Bob | Sony | 2014-04-20 | Buy |
| … | … | … | … |

- **Alibaba organizes a competition to look for even better prediction models**

- **Given the log of user behaviors (including CLICK, BUY, BOOKMARK, and ADD-TO-CART) of a certain period**

- **Predict which users will buy which brands at a later time**

**We are the No.1 out of 7000+ teams in season 1.**

# Smart Tour in Canton Tower

A mobile phone app with augmented reality

**Your smartphone will accompany you go shopping**

- **Indoor navigation**
- **Personalized and location-based recommendation**
- **Finding your cars in big parking lots**

**Better shopping experience with very low cost!**

- It is estimated that there are **more than 2500** false insurance claims each year in UK.

- **Each** false claim could cause a loss of **up to $300,000 HKD**.

- In the market, there are **old and expensive** anti-fraud solutions by FICO, SPSS, etc.

- Our solution uses **new machine learning technologies**.

**Our solution detects more frauds while cuts the cost by more than a half.**
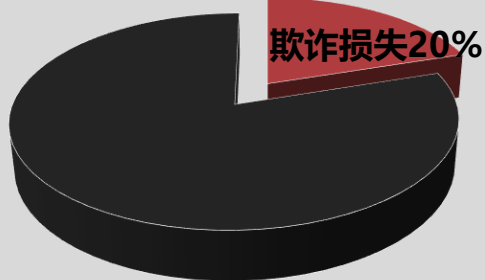
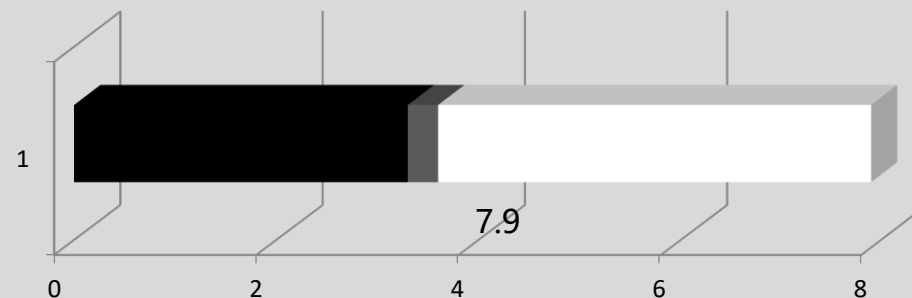# Fraud Detection for Automobile Insurance

| 2013 年保费收入居于前 5 位的商业保险险种 | | | | | | |
|---|---|---|---|---|---|---|
| 保费排名 | 险种 | 保费收入 | 保险金额 | 赔款支出 | 未到期责任准备金 | 未决赔款准备金 | 分险种利润表承保利润 |
| 1 | 机动车辆保险 | 350,583.42 | 38,410,740.36 | 196,578.28 | 145,825.68 | 108,060.02 | -3,285.97 |
| 2 | 企业财产保险 | 49,429.71 | 119,063,252.44 | 31,968.74 | 16,444.05 | 299,200.00 | -5,059.11 |
| 3 | 责任险 | 48,294.09 | 21,807,720.77 | 12,459.24 | 15,574.38 | 35,403.74 | 2,326.12 |
| 4 | 意外险 | 39,322.31 | 5,164,595,780.71 | 8,713.03 | 8,384.49 | 8,363.93 | -2,101.88 |
| 5 | 货运险 | 35,933.58 | 84,998,679.25 | 13,896.63 | 2,737.79 | 17,618.38 | 3,577.47 |

单位：万元

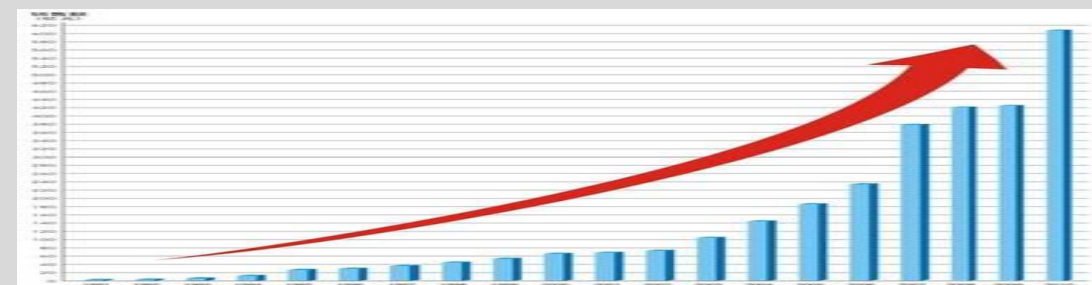An auto insurance company had compensation expenses 1.97B, with net loss 33M in 2013

欺诈损失20%

CIRC statistics show that about 20% of claims are fraudulent

Save up to 78M from the loss of 390M fraud

With 3M system cost, the company made profit of 42M

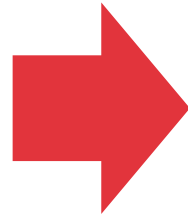Fraud detection rate increased by 20%, you can change from a loss of 9.4 ‰ to 12.2 ‰ profit
Benefitted from Big Data, the company becomes profitable!

9

**User Data**

Age、Sex、Job
Type of Phone
Phone Log
Traffic Statistics
Complaint Record
Home Zone
Location
Network Time
Payment Record
……

Customer Loyalty Prediction Model

Personalized Customer Service

In Q1 2011, customer drop out rate reduced by 50% in the US

Mobile Internet allows real-time communication with the Internet

- Always connected

- Greater user stickiness

- Longer access time

- Lower cost of participation

Significantly change to user behavior

FB: $16B          >$10B



**New IT Business**



**Smart Shopping Mall**



**Microfinance Companies**



**Data Center Infrastructure**



**Data-based Precise Marketing**



**Credit Scoring Service**

**12**

# It requires new information technologies



**Big Data Collection**



**Large-scale Data Mining**



**Advanced Machine Learning**



**Data** → **Knowledge** → **Insight**

# The point of Big Data is to make sense of it

# Challenges of Harnessing Big Data for Business Value

## Data Collection

- Extracting siloed data
- **Understanding metadata**
- Mixed structured and unstructured data
- Data cleaning
- **Data calibration**
- Data integration
- …

## Data Processing

- Data loading
- Building large indexes
- **Parallel algorithms design**
- **Fault-tolerance**
- Adapt to new hardware
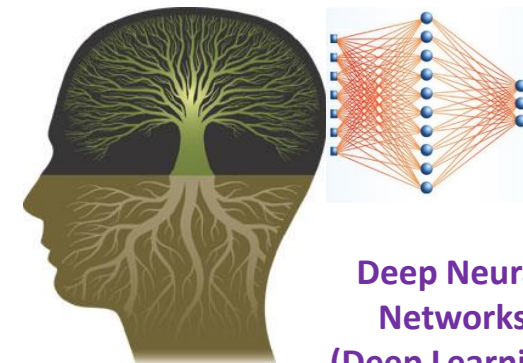- Data compression
- **Real-time response**
- …

## Data Mining

- Statistical analysis
- Data clustering
- **Predictive modeling**
- Ensemble of models
- Abnormal detection
- **Unsupervised learning**
- Data visualization
- …

## Winners in the Big Data era would be those who:

- Have as much as possible data
- Have a creative mind of what the data can do (the value of the data)
- Know how to extract knowledge and gain insight from data
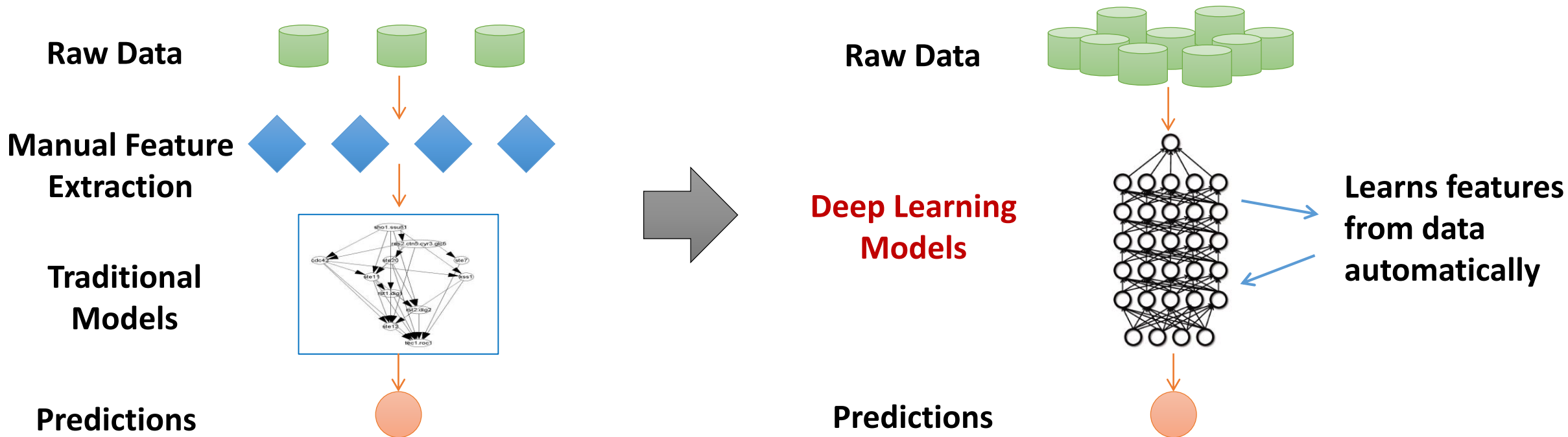- Have the computation resources and know how to process huge amounts of data



Deep Neural Networks (Deep Learning)

Raw Data

Manual Feature Extraction

Traditional Models

Predictions

Deep Learning Models

Raw Data

Learns features from data automatically

Predictions

**Amazing accuracy (> 85%) for image recognition; yet many works remain to be done to apply it for business data.**

## Models are no secret
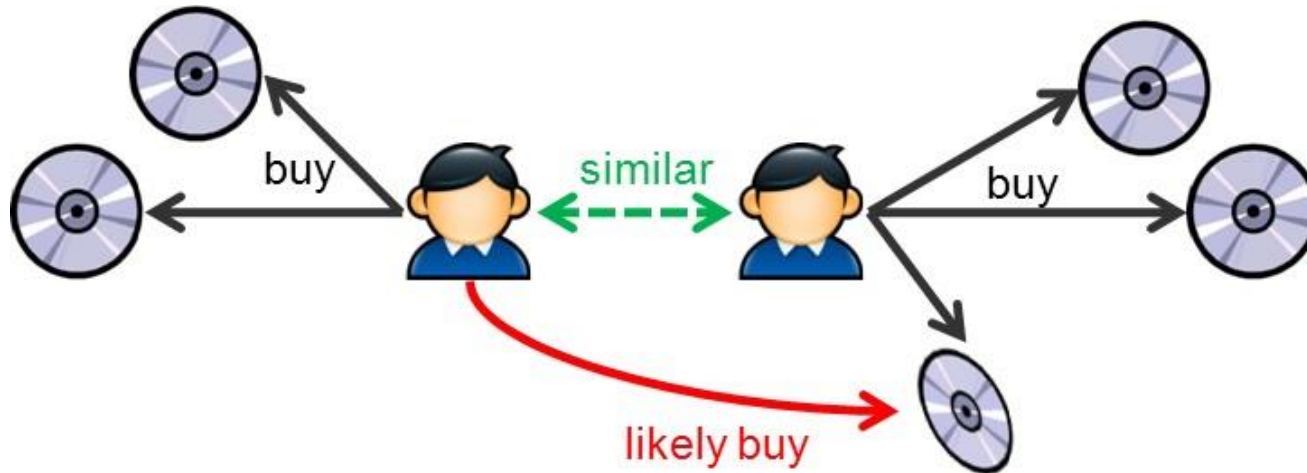
- They are all textbook knowledge.
- Alibaba even let the players know what models they use.

## The point is how to well use them

- How similar is similar?
- How many similar users/brands to use?
- How to tune the parameters?
- …



**Rocket science? No! It's all art!**

## To collect more data

- "Data will become more valuable than you thought when you collect it."

  — *Dr. Jian Wang, CTO of Alibaba Corp.*

## To profit from data

- Treating big data as profitable assets
- Figuring out how data will help the business

## To hire experienced data scientists

- Data science is more of an art than a science.



**Harvard Business Review**

THE MAGAZINE    BLOGS    VIDEO    BOOKS    CASES    WEBINARS    COU

**Data Scientist: The Sexiest Job of the 21st Century**

by Thomas H. Davenport and D.J. Patil

Comments (87)

Newsroom \ Announcements \ Gartner Says Big Data Creates Big Jobs: 4.4 Million IT Jobs Globally...

Press Release    Share:    Like 38    Tweet 36    Share 36    +4

ORLANDO, Fla., October 22, 2012    View All Press Releases ›

**Gartner Says Big Data Creates Big Jobs: 4.4 Million IT Jobs Globally to Support Big Data By 2015**

Analysts Discuss Key Issues Facing the IT Industry During Gartner Symposium/ITxpo 2012, October 21-25, in Orlando

## Build data center

- Not affordable to SMEs.

## Collect more data

- Collect and publish public data.

## Train more data scientists

- Preferential policies for Big Data initiatives
- More research funding in Big Data

Flow of people
(e.g., octopus data)

Taxi and bus GPS data

Weather and
environment data

Mobile phone
base station data

**Collecting and
providing data as a
public service**

City data centers

**Taking care of
privacy issues**
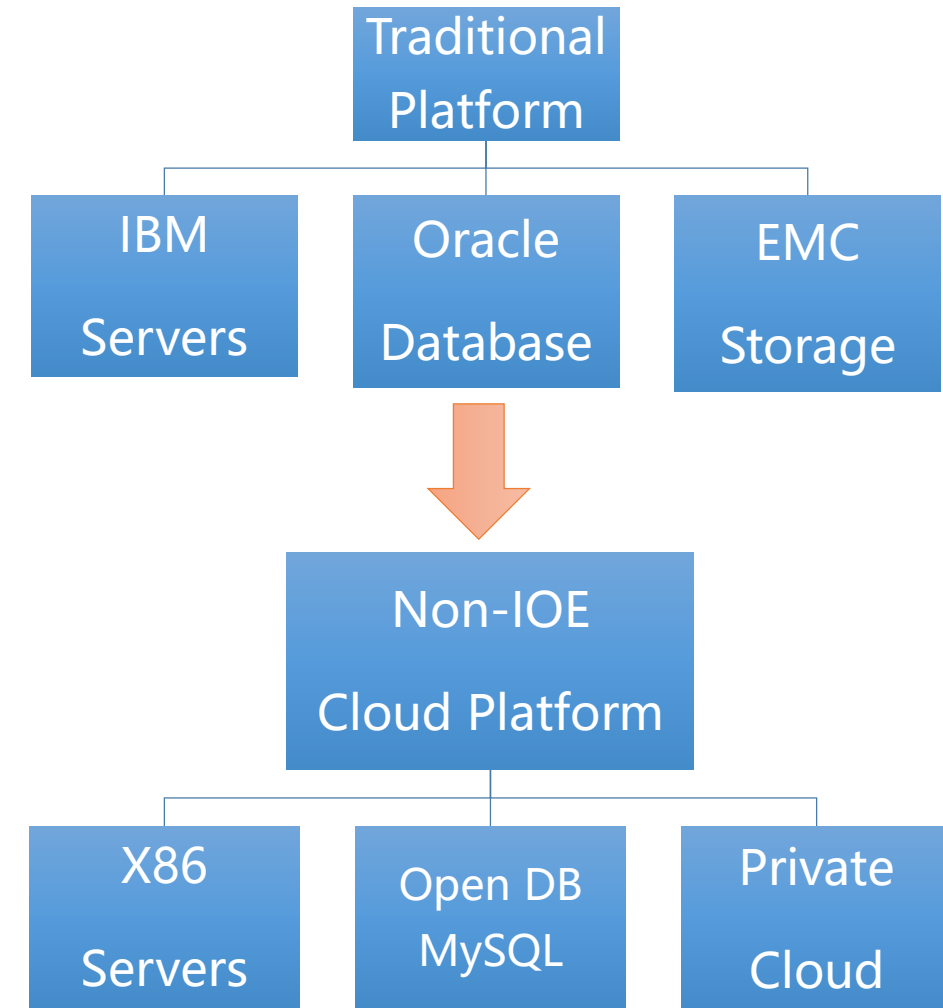
Business

Research

Public administration

Personal use

**IBM servers, Oracle database, EMC storage systems**.
IOE have dominated the hardware and software of enterprises for data management, especially those big enterprises, financial companies in China.

Traditional Platform

IBM Servers

Oracle Database

EMC Storage

Non-IOE Cloud Platform

X86 Servers

Open DB MySQL

Private Cloud

# Data-intensive vs. Computation-intensive

**In the past ten years, Big Data is mainly about data-intensive processing.**

**In the next ten years, Big Data may be more about computation-intensive processing.**

Easy tasks, e.g., simple queries, building inverted index

Difficult tasks, e.g., training complex machine learning models

Petabytes of (text) data

Relatively small amount of (feature) data extracted from raw data

Thousands of commodity (cheap) servers with lots of hard disks

Hundreds of powerful servers with lots of processing units/cores

Example: Google's cluster for web search
More than 10,000 low-end servers

Example: IBM Watson
Less than 100 high-end servers

23

- Data are always biased, regardless of its size.
- In most cases, big data supports decision making; yet it is unwise or even dangerous to let data alone make decisions for us (human beings).
  - Correlation rather than causality
- Value big data, do not ignore "small data".


Big Data: A young man, a conqueror, and a reaper, who still needs practicing, developing, receiving and understanding

- New technologies do bring other concerns, such as privacy. We should face them and solve them, other than escape from them.

- Technologies will move the world forward, and there is no way back.

- The history has told us that those companies who overlook the technology advances will be out of business, faster than you thought.



Do You Still Remember Kodak?

# Big Data Will Change Business

**In the near future:**

**Big Data Should Be and Could Be
A Core Competency of Your Enterprise**

**Thank You !**