

# 大數據的解決方案及其應用

雲計算在統計及Data Mining研究

應用及前沿綜述

雲計算下的R前沿探討 - 兼談 數據分析的未來

## 謝邦昌

輔仁大學商學研究所博士班 所長

統計資訊學系教授

中華資料採礦協會 榮譽理事長

2014年10月02日

[025674@mails.fju.edu.tw](mailto:025674@mails.fju.edu.tw)

[WWW.CDMS.ORG.TW](http://WWW.CDMS.ORG.TW)



刁院士：統計與數學要相互欣賞

統計要跟各個領域  
做朋友



統計趨勢 Statistics Trend 趨勢統計 Trend Statistics

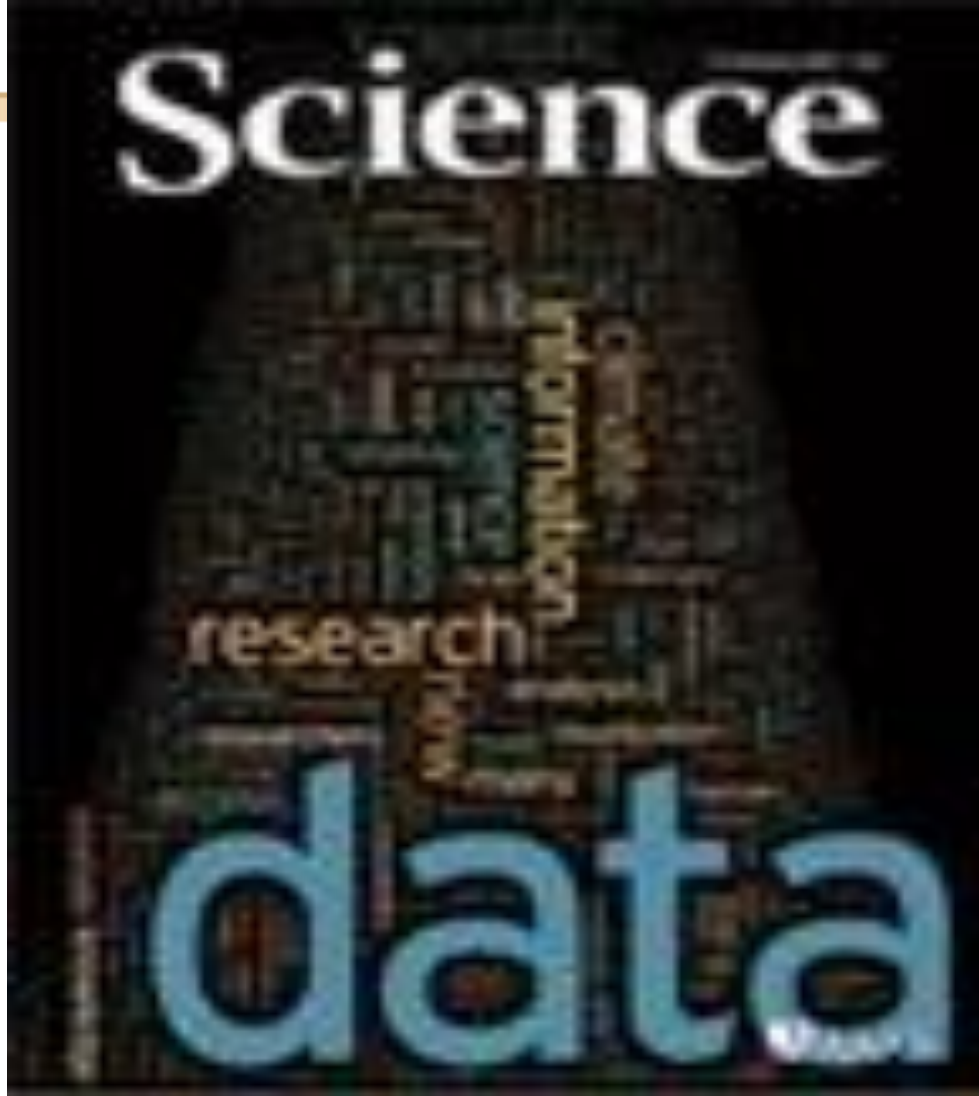
數據分析科學的過去、現在、未來

--統計是數據科學

Data Science

Dataology





# Science

-- ( Dealing with DATA ) 11 FEBRUARY  
2011 VOL 331 SCIENCE [www.sciencemag.org](http://www.sciencemag.org)



Revolution  
in Egypt

**Joe Klein:** What the U.S. should do  
**On the Street:** Hope meets anxiety  
**Muslim Brotherhood:** What it wants

**Oscars:**  
Portraits of  
star power

TIME

# 2045

The Year Man Becomes Immortal\*

BY LEV GROSSMAN

\* if you believe  
humans and  
machines will  
become one.  
Welcome to  
the Singularity  
movement



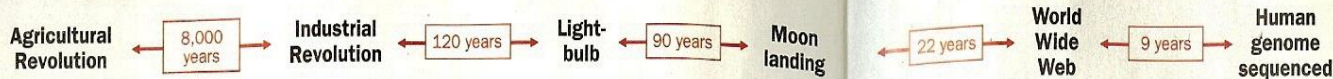
9 771064 030005

www.time.com

BRUNEI ..... B\$ 9.00 HONG KONG ..... HK\$ 95.00 KOREA ..... ₩ 7,900.00 MONSIEUR .....  
CAMBODIA ..... US\$5.00 INDIA ..... ₹ 110.00 MACAU ..... MOP 95.00 MALTA ..... M€ 5.00  
CHINA ..... RMB 90.00 INDONESIA ..... Rp 95,000 MALAYSIA ..... RM 11.00 NEPA .....  
US\$ 10.00 PHILIPPINES ..... ₱ 550.00 SINGAPORE ..... S\$ 10.00 THAILAND ..... ฿ 55.00  
US\$ 10.00 U.S. MAIL PERMIT NO. 1000 NEW YORK, NY 10108



# 1 The accelerating pace of change ...

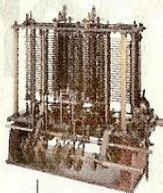


# 2 ... and exponential growth in computing power ...

Computer technology, shown here climbing dramatically by powers of 10, is now progressing more each hour than it did in its entire first 90 years

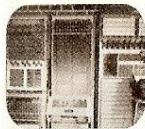
## COMPUTER RANKINGS

By calculations per second per \$1,000



### Analytical engine

Never fully built, Charles Babbage's invention was designed to solve computational and logical problems



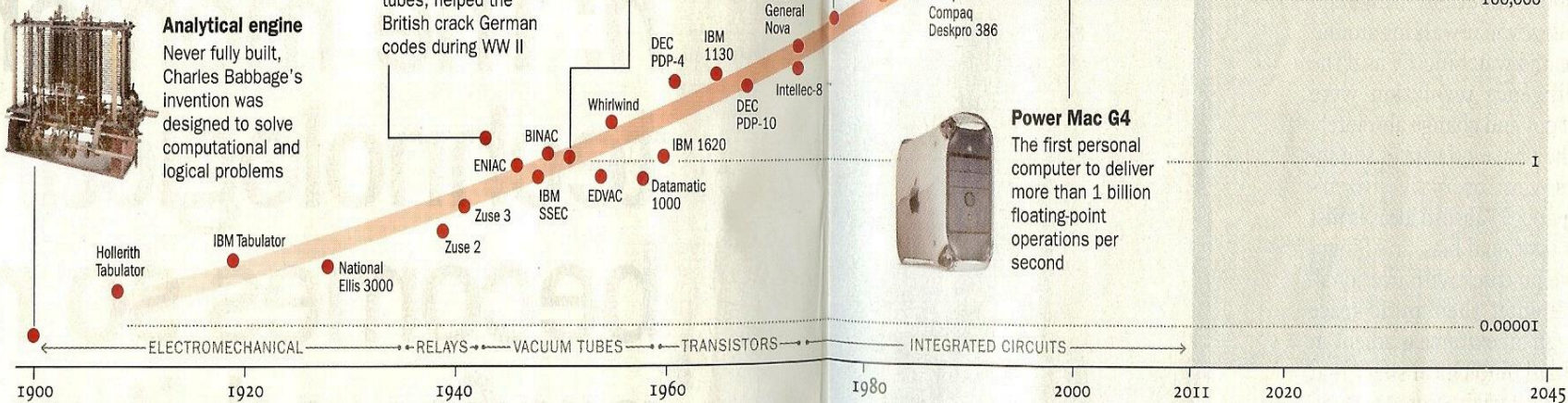
### Colossus

The electronic computer, with 1,500 vacuum tubes, helped the British crack German codes during WW II



### UNIVAC I

The first commercially marketed computer, used to tabulate the U.S. Census, occupied 27 cu m



# 3 ... will lead to the Singularity



**Apple II**  
At a price of \$1,298, the compact machine was one of the first massively popular personal computers



**Power Mac G4**  
The first personal computer to deliver more than 1 billion floating-point operations per second

**2045**  
Surpasses brainpower equivalent to that of all human brains combined

Surpasses brainpower of human in 2023



Surpasses brainpower of mouse in 2015

on, there's no reason to think computers

Probably. It's impossible to predict the

idea; it's a serious hypothesis about the

he called an "intelligence explosion":





「電腦vs. 人腦」益智遊戲大賽經過三天激戰後，IBM的超級電腦華生（Watson，圖中）最終擊敗人類，獲頒100萬美元獎金。IBM將把這筆獎金捐給世界展望會等慈善機構。

（美聯社）

# 雲端運算的演化

**Super  
Computer**

**Cluster  
Computing**

**Distributed  
Computing**

**Grid  
Computing**

**Utility  
Computing**

**Cloud  
Computing**

## 雲端運算

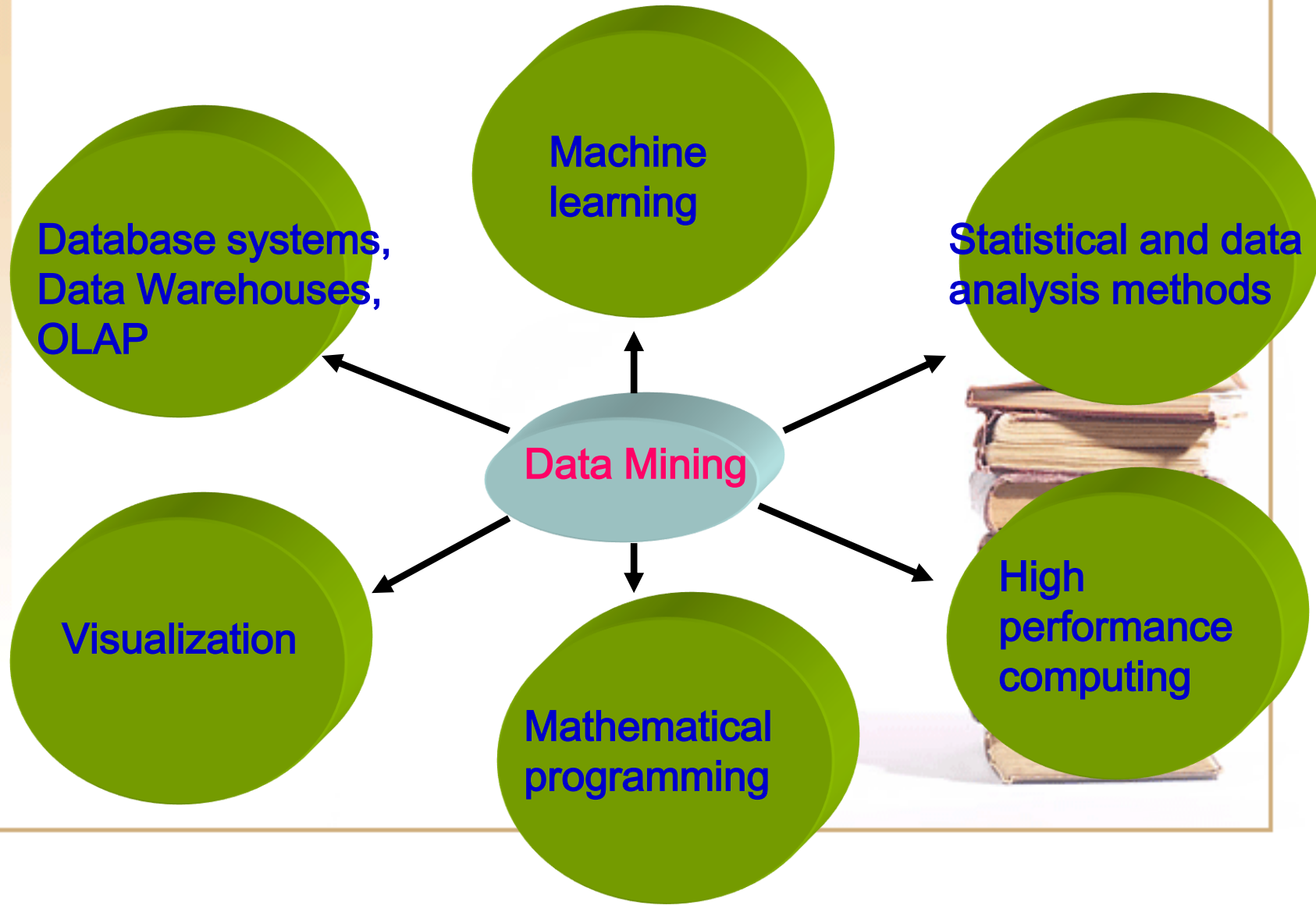
- 透過網路將龐大的運算處理常式自動分拆成無數個較小的副程式，再交由多部伺服器所組成的龐大系統經搜尋、運算分析之後將處理結果回傳給用戶
- 雲 ~ = 網路
- Google: MapReduce、GFS及BigTable



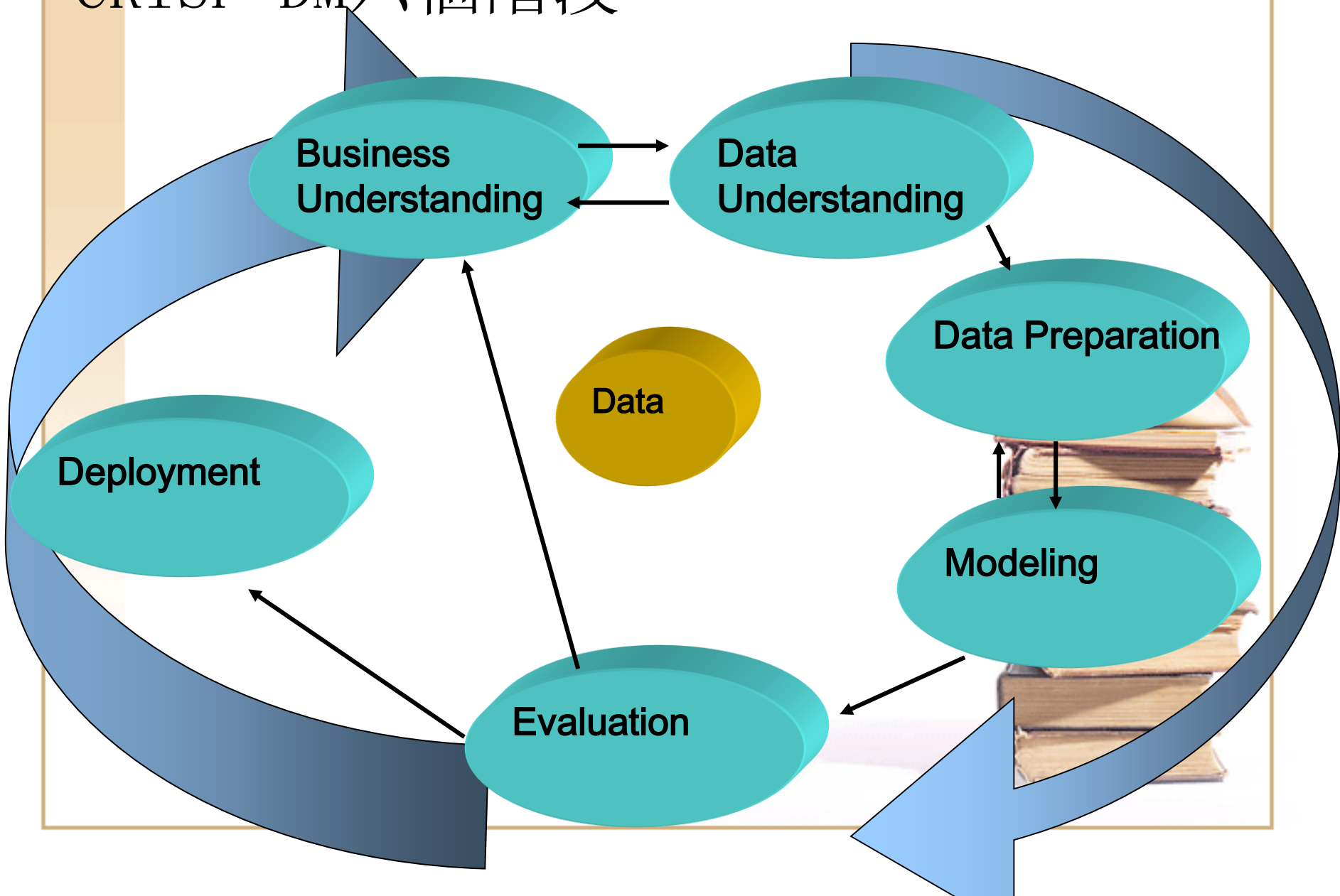
# 比較表

服務 屬性	Amazon EC2	Google App Engine	Microsoft Azure	Yahoo Hadoop
架構	Iaas/Paas	Paas	Paas	Software
服務型態	Compute/ Storage	Web application	Web and non-web	Software
管理技術	OS on Xen hypervisor	Application container	OS through Fabric controller	Map / Reduce Architecture
使用者介 面	EC2 Command- line tools	Web-based Administratio n console	Windows Azure portal	Command line and web
APIs	yes	yes	yes	yes
收費	yes	maybe	yes	no
程式語言	AMI(Amazon Machine Image)	Python	.NET framework	Java,

# Data Mining包含六大領域

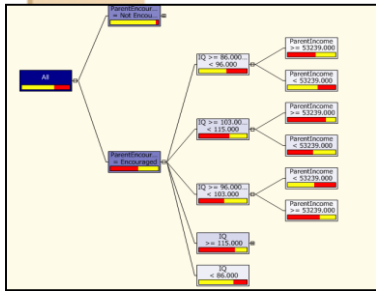


# CRISP-DM六個階段



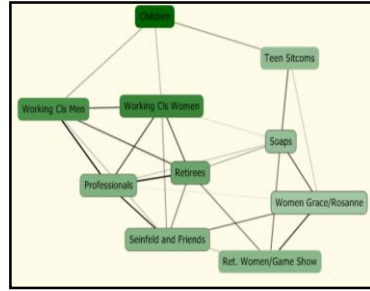


# 完整的演算法 SQL Server 2008已提供 SQL Server 2014在雲端

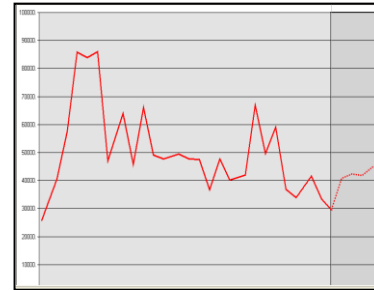


決策樹

SQL Server 2000已提供



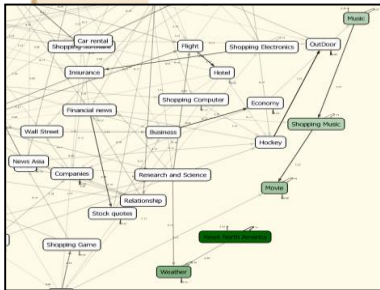
群集



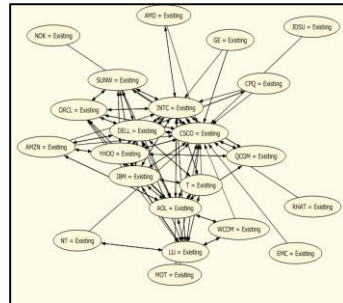
時間序列

Attributes	Values	Favors Professional/Techn.	Favors Service Workers
Education Years	15-20	█	
Education Years	12-13		█
Education Years	7-12		█
netion hbQJUNG AND THE RES.	Missing	█	
netion hbQJUNG AND THE RES.	Existing		█
netion hbQS THE WORLD TURN.	Existing		█
netion hbQS THE WORLD TURN.	Missing		█

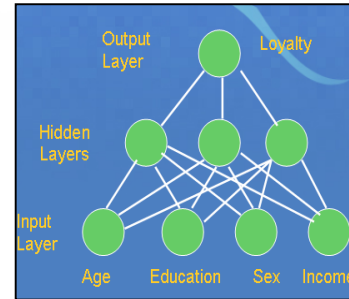
貝氏機率分類



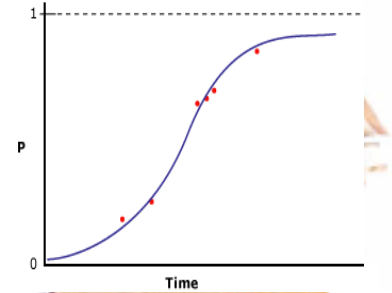
時序群集



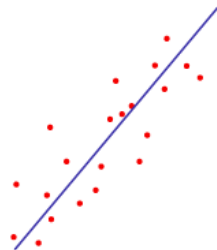
關聯規則



類神經網路



羅吉斯回歸



線性回歸

When upgrading to Microsoft® SQL Server™ 2000, you can upgrade servers in your org at a time; however, when servers are used for **publishing**, you must upgrade the Distrib Publisher second, and then Subscribers. Upgrading servers one at a time following this recommended when a large number of Publishers and Subscribers exist because you can **publish** data even though servers are running different versions of SQL Server. You can publications and subscriptions with servers running instances of SQL Server 2000, and at subscriptions created in SQL Server 6.5 or SQL Server 7.0.

When using transactional **publishing**, you can upgrade Subscribers before the Publisher, using immediate updating with snapshot **publishing** or transactional **publishing**, there are upgrade recommendations in this topic under Upgrading and Immediate Updating.

You can upgrade **publishing** servers running SQL Server 6.5 or SQL Server 7.0 to SQL 5 the server is running SQL Server 6.5, you do not need to upgrade it to SQL Server 7.0 before upgrading to SQL Server 2000.

**IMPORTANT** When upgrading servers configured for **publishing** to SQL Server 2000, if compatibility level must be set to 70 (version 7.0 compatibility) or later. If you have a running in 65 (version 6.5) or an earlier compatibility level, temporarily change them during the upgrade process.

When the Publisher or Subscriber is running in 65 or an earlier compatibility level due to SQL Server 2000, error 15048 will be raised stating that the operation is supported Server version 7.0 or SQL Server 2000.

For more information about setting the backward compatibility level, see [SQL Server 200 Server version 6.5](#).

If you are upgrading **publishing** on a failover cluster, you must uncluster the previous one before upgrading. Unclustering the previous installation means that you must delete all **publishing**, and reconfigure it after upgrading to SQL Server 2000. This will not requirement when upgrading SQL Server 2000 to future releases.



# 常用的Data Mining及統計學習方法-1

Binary Classifier (二元分類)

Numeric Predictor (數值預測)

Time Series (時間序列)

C&R TREE (分類回歸樹)

Quick Unbiased Efficient Statistical Tree (QUEST決策樹模型)

CHAID (分類樹)

Decision List (決策樹列表)

Regression (線性回歸分析)

PCA/Factor (主成分分析)

Neural Net (類神經網路)

C5.0 (決策樹)

Feature Selection (特徵選取)

Discriminant Analysis (判別分析)

Logistic (羅吉斯回歸)

Generalize Linear Model (廣義線性模型)

Cox Regression



# 常用的Data Mining及統計學習方法-2

Support Vector Machine (SVM支持向量機)

Bayes Net (貝氏分類器)

SLRM (自我學習反應模型)

GRI關聯

Apriori關聯

CARMA關聯(連續交易)

Sequence Clusterc序列關聯

K-Means (K-Means分群)

Kohonen (自我組織化)

Two-Step (二階段)

Anomaly (異常檢測)

Random Forests (隨機森林)

ICA (獨立成分分析)

Multivariate adaptive regression spline (MARS多元適應性回歸平滑)

Pmml(預測模型標記語言)

Boosting





# 使用軟體 常用

SQL server 2014

SPSS 22 (PAWS) --IBM

SAS

SQL 2014+Excel (2013)-Data Mining

Add-in

Clementine 12.0

Statistica 14.0

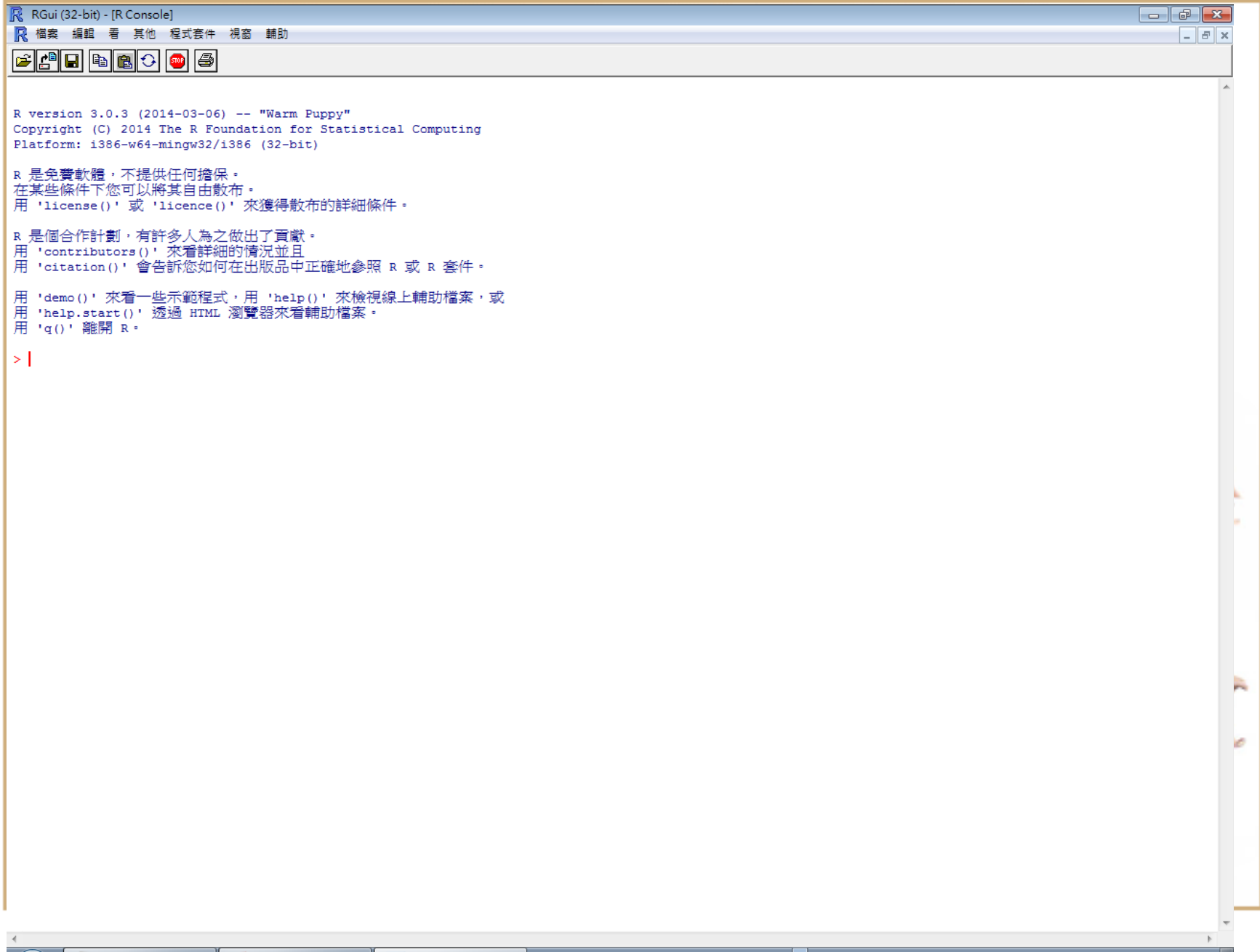
WEKA

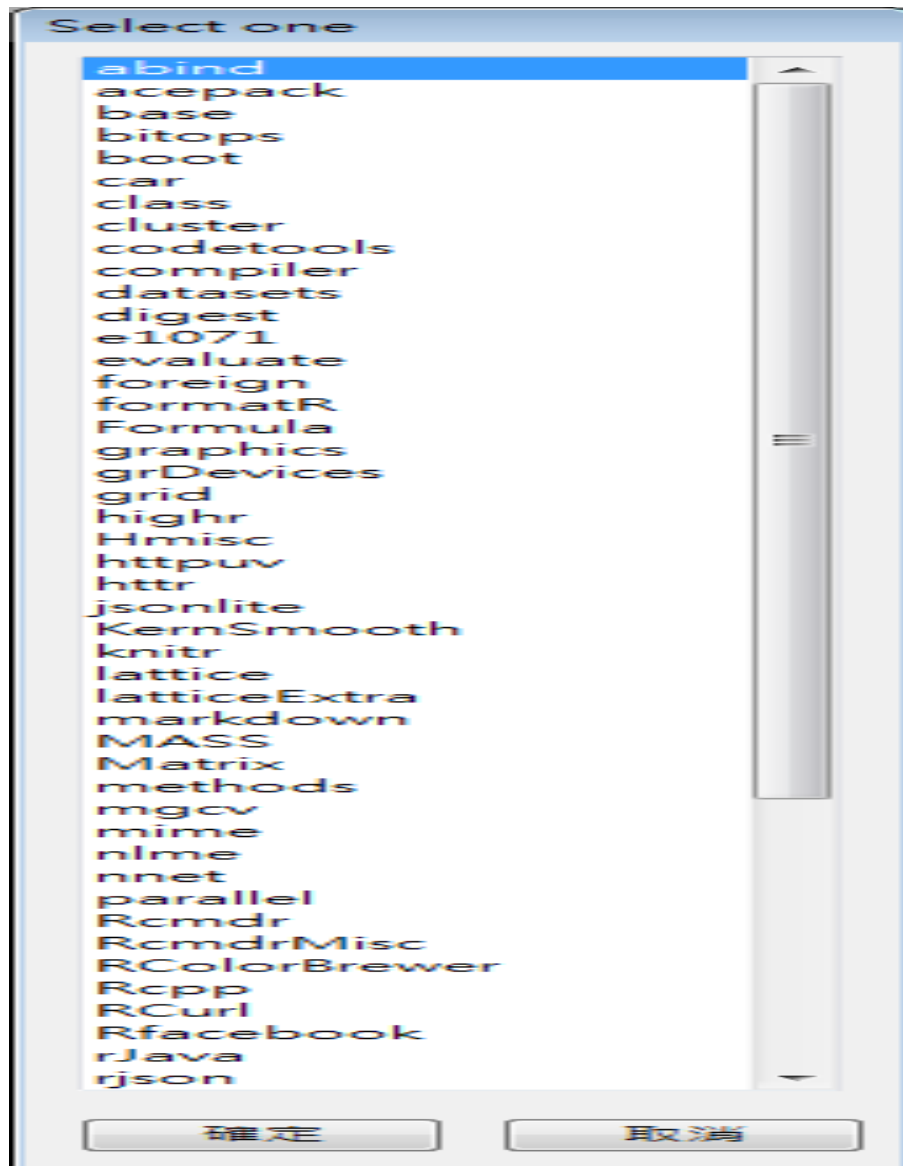
**R → Cloud R**

**R+Excel ADD-IN** ..... 還有更多雲端軟

**體**









# R Excel

## 數據採礦發展趨勢



波音 DC-8 客机

- RExcel

- RExcel之創始

- RExcel之啟動

- RExcel之應用

- 數據導入

- 數據分析

- 結果保存



# statconn之“幕后黑手”

## (The masterminds behind statconn)



- Thomas Baier (1971-)
- 在不同環境中應用R
  - R/Scilab (D)COM Server
  - RExcel (1998)



- Erich Neuwirth (1948-)
- RExcel的主要作者



University of  
Vienna RExcel之創始

• <http://rcom.univie.ac.at/>

# 數據採礦在各產業的應用

- **金融服務業**

客戶貢獻度分析、信用評分、風險評估、客戶區隔、交叉行銷等。

- **保險業**

顧客貢獻度分析、信用評分、風險評估、客戶區隔、交叉行銷、客戶流失分析和詐欺偵測等。

- **電信業**

顧客貢獻度分析、信用評分、客戶區隔、交叉行銷、客戶流失分析、銷售預測和詐欺偵測等。





# 數據採礦在各產業的應用

- **製造業**

客戶貢獻度分析、品質管制、行銷績效分析、生產分析和存貨分析等。

- **零售業**

客戶忠誠度、客戶區隔、購物籃分析、定價分析、交叉行銷和銷售預測等。

- **生物科技、醫療保健、航太空業、環境、法律等**



# 商業智慧的核心

- 如何收集數據
  - 營運數據，市場調查數據，固定 **Panel** 追蹤
- 如何管理數據
  - ETL，Data warehousing
- 如何從數據中獲取智慧
  - Data Mining，OLAP，Statistics
- 如何應用智慧
  - 行銷策略，主管決策，互動化 **CRM** 機制



- 雲端運算可以實現適應端通過在線上傳數據或購買數據，通過雲數據倉庫，進行數據倉庫建模或數據抽取，線上支付使用數據採礦工具和商業智慧相關處理軟體



# IIaaS是SaaS的延伸 $I^2aaS$

- 數據採礦和商業智慧的原理相似，均由數據提供資訊、產生知識，再由知識累積智慧。而雲端運算可以使這個過程在網際網路上得以實現。也就是說雲端運算可以提供基於SaaS的知識與智慧分析的服務

(Information&Intelligence as a Service) ,  
簡稱IIaaS ;  $I^2aaS$ ，它是SaaS  
的延伸。





# 雲端運算產業類型

IIaaS  $I^2$ aaS

Information & intelligence as a Service

SaaS

Software as a Service

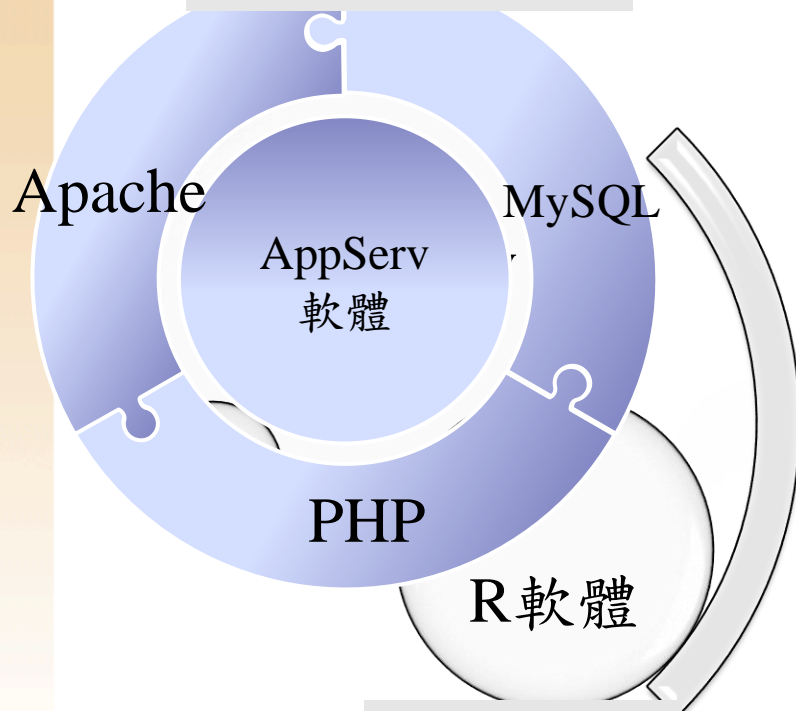
PaaS

Platform as a Service

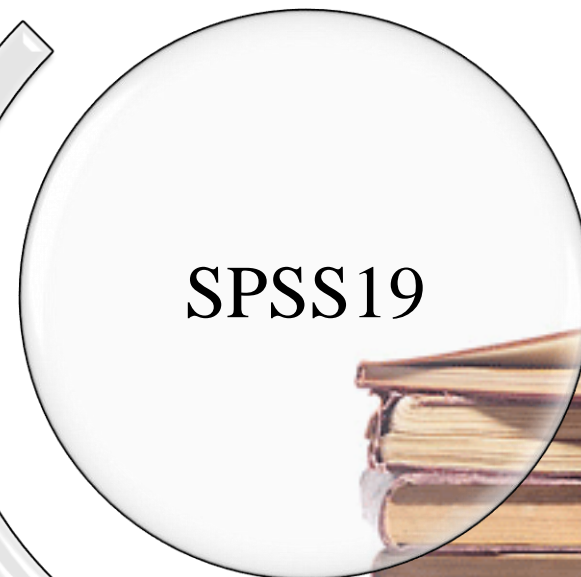
IaaS

Infrastructure as a Service

## 系統環境建置



## 統計分析



智慧型稅務選案系統

數據整理

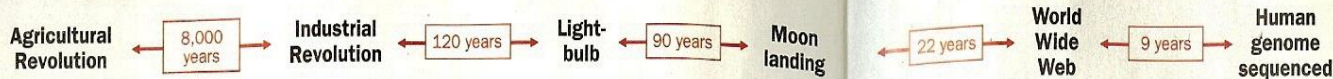
# 科技「瘋」雲，再輝煌十年

- 「雲端運算 (Cloud Computing) 即將引爆商業革命，改寫遊戲規則，」2009年6月，美國《BusinessWeek》如此寫道。「雲端運算讓企業節流，也可以變得有創意，」2009年7月，美國《哈佛商業評論》製作後風暴時代首部曲專題，撰文分析。「雲端運算將是一朵長長的雨雲 (nimbus)，讓企業更靈活，」從去年10月以來，英國《經濟學人》也陸續有著相關報導。近期以來，關於雲端運算議題不斷發燒，不僅全球重要媒體關注，各國企業更是積極投入。





# 1 The accelerating pace of change ...

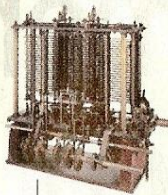


# 2 ... and exponential growth in computing power ...

Computer technology, shown here climbing dramatically by powers of 10, is now progressing more each hour than it did in its entire first 90 years

## COMPUTER RANKINGS

By calculations per second per \$1,000



### Analytical engine

Never fully built, Charles Babbage's invention was designed to solve computational and logical problems



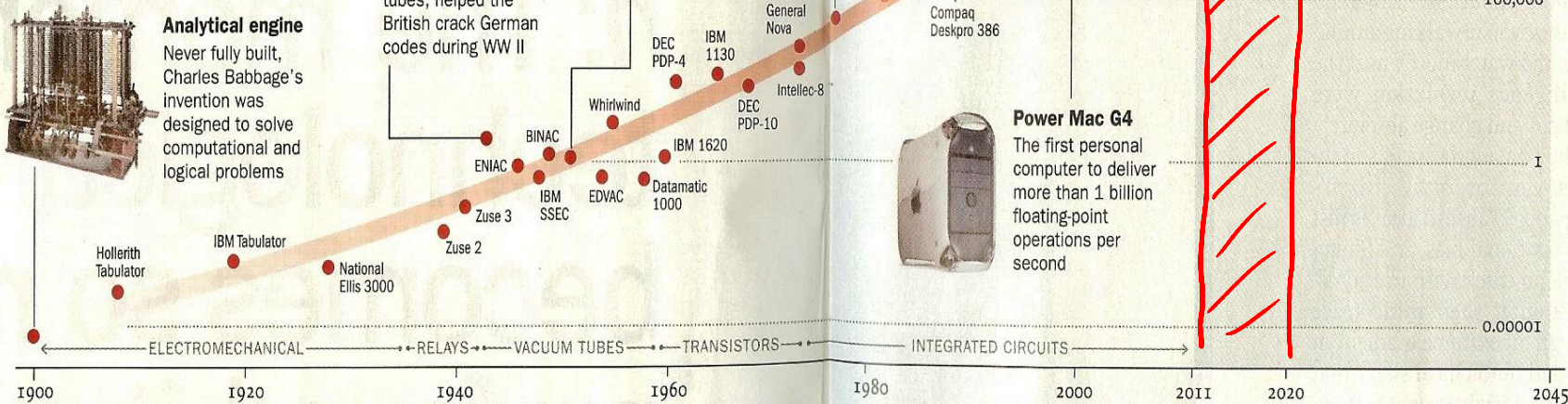
### Colossus

The electronic computer, with 1,500 vacuum tubes, helped the British crack German codes during WW II



### UNIVAC I

The first commercially marketed computer, used to tabulate the U.S. Census, occupied 27 cu m



# 3 ... will lead to the Singularity



**Apple II**  
At a price of \$1,298, the compact machine was one of the first massively popular personal computers



**Power Mac G4**  
The first personal computer to deliver more than 1 billion floating-point operations per second

on, there's no reason to think computers Probably. It's impossible to predict the idea; it's a serious hypothesis about the he called an "intelligence explosion":



漫步雲端，  
任重而道遠！

只在此山中 雲深不知處

