

# 第一讲

## 机器学习简介

**夏睿**

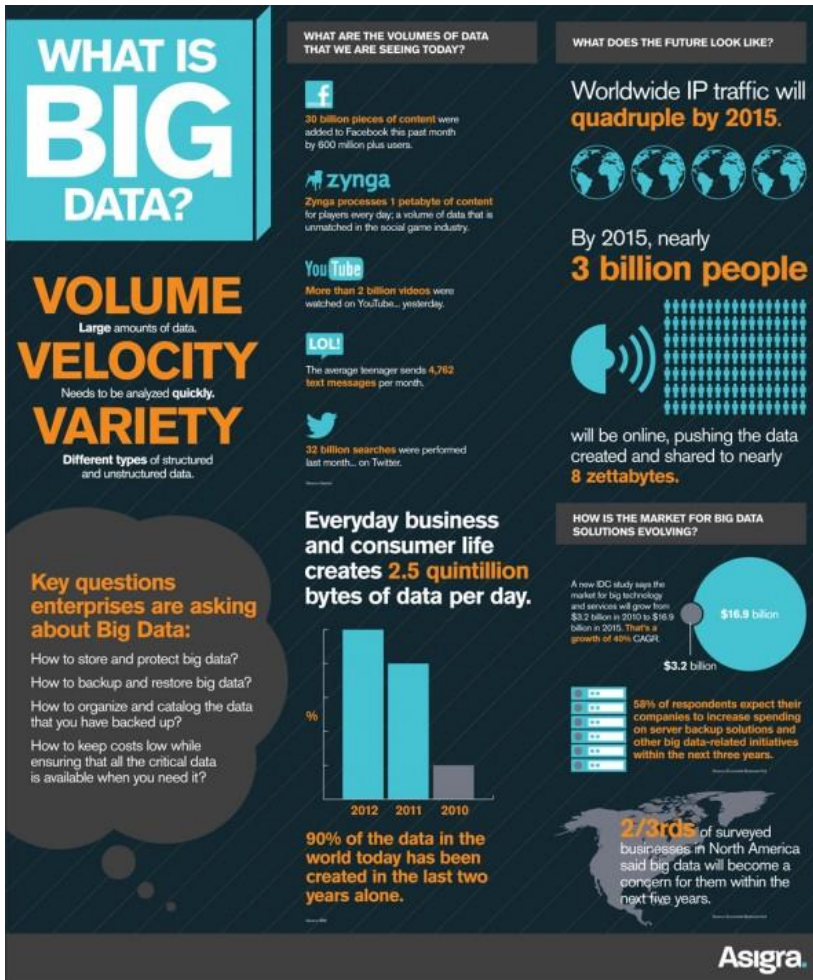
计算机科学与工程学院

南京理工大学

[rxia@njust.edu.cn](mailto:rxia@njust.edu.cn)

<http://www.nustm.cn/member/rxia>

# 大数据时代



- 每天大约有50亿tweets被传送，也就是说每秒钟有超过5700条tweets。
- Facebook上有超过11.5亿的活跃用户在互动。
- 超过50亿人群正在打电话、发短信、发微博、在手机上浏览网站。
- VISA每天大约有172,800,000的信用卡业务需要处理。
- United Parcel Service（美国联合包裹服务）每天能收到大约395,000,00包裹。
- RFID（射频识别系统产生高达常规条形码系统数据的1,000倍。

# 什么是机器学习?



数据



机器学习



认知

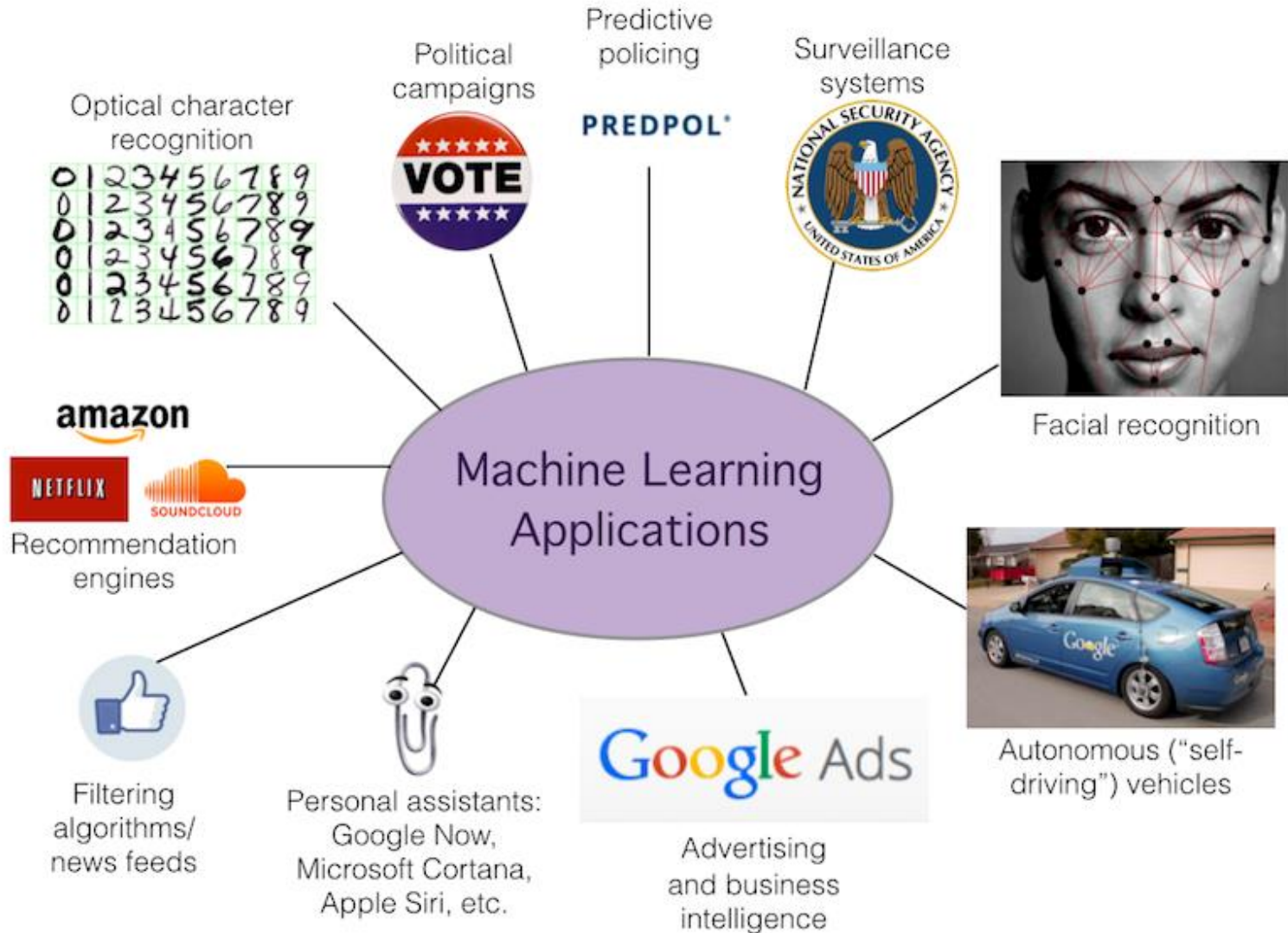
# 机器学习的定义

- Arthur Samuel (1959)定义机器学习是一个“无需明确编程而给予计算机学习能力的研究领域。”
- Tom M. Mitchell (1997)提供了一个更广泛的引证，更规范的定义是：  
“如果它在任务T中的性能被P测试，随着经验E的增加而提高，则被认为是计算机程序是从经验E中学习某些类别的任务T和性能度量P。”



机器（电脑）了解我们（人类）可以学习什么（从数据）？

# 机器学习的应用



# 文本分类



U.S. **Top Stories**

- Top Stories
- Starred ☆
- World
- U.S.
- Business
- Sci/Tech
- Entertainment
- Sports
- Health
- Spotlight
- Most Popular


► All news  
[Headlines](#)

★ **Obama's budget proposal draws rapid fire from legislators**  
USA - [Richard Wolf](#), [Steve Hebert](#) - 2 hours ago  
WASHINGTON - President Obama's proposed \$3.8 trillion budget ran into immediate trouble in Congress on Monday among lawmakers who said it tries to do too much while cutting the deficit too little.  
[Video: Obama Budget Would Create Highest-ever Deficit](#) The Associated Press  
[Wealthy Face Tax Increase](#) Wall Street Journal  
[Milwaukee Journal Sentinel](#) - [Bradenton Herald](#) - [ABC News](#) - [Daily Caller](#)  
[all 4,573 news articles >](#) [Email this story](#)

☆ **Md. stands to gain despite Obama budget cuts**  
Baltimore Sun - [Paul West](#) - 1 hour ago  
WASHINGTON - President Barack Obama wants to end the nation's troubled program to return a but NASA officials indicated Monday



The Hindu



Google in:spam

Gmail -

COMPOSE

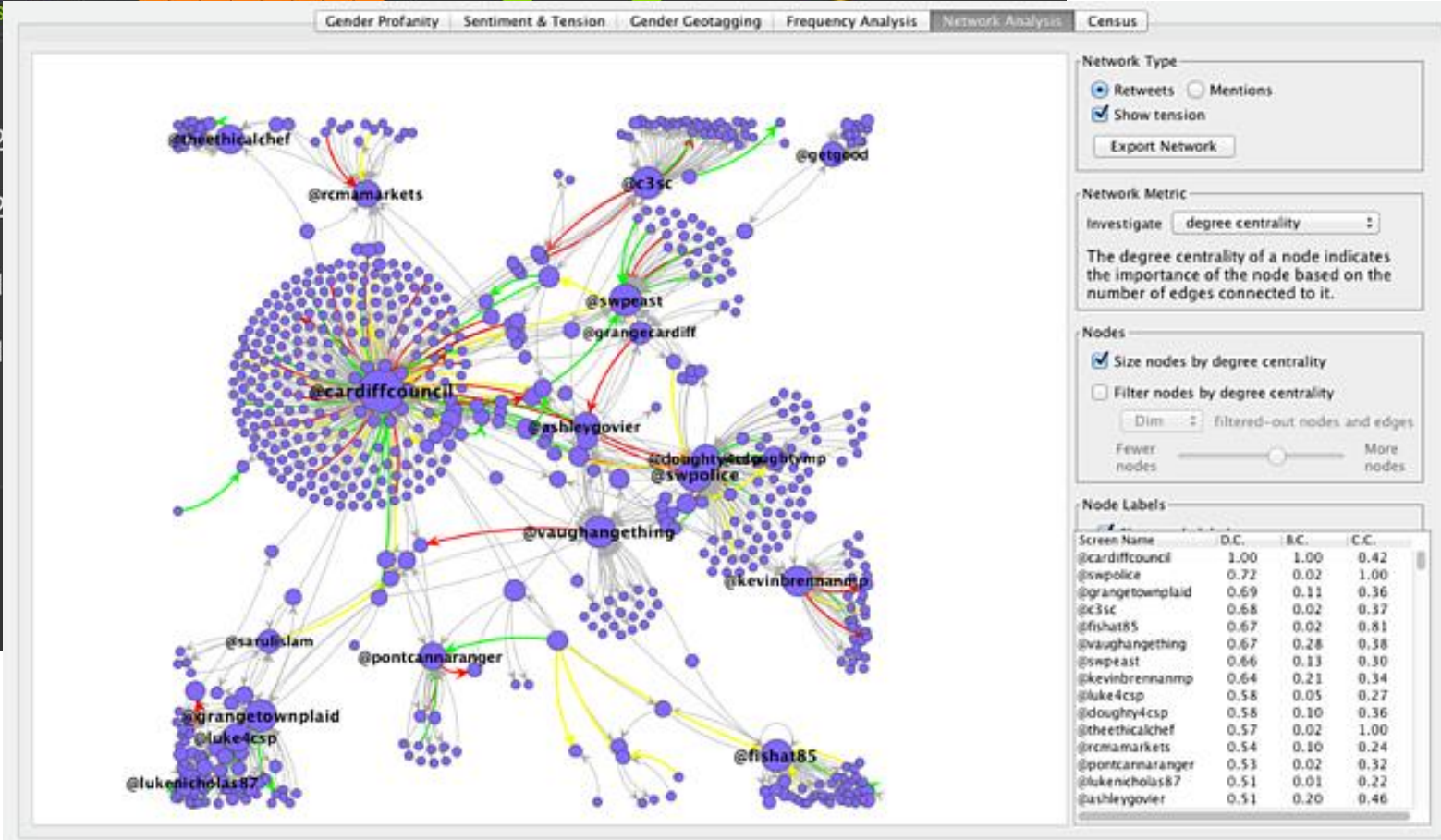
- Inbox (7)
- Starred
- Important
- Sent Mail
- Drafts (15)
- All Mail
- Spam (46)
- Trash
- Circles

<input type="checkbox"/>	☆	me	Delete all spam messages now (messages that have been in Spa
<input type="checkbox"/>	☆	no1.gr	New submission from Quick Poll: Facebook Pre-Fill - I would u
<input type="checkbox"/>	☆	PayPal	Προσπαθήστε το κινητό σας... - Ean den mporcite na delte to ne
<input type="checkbox"/>	☆	EdFed	Your PayPal account has been limited! - Warning Notification De
<input type="checkbox"/>	☆	LoopGalaxy	"What NOT TO DO During Your Interview" - To ensure prompt d
<input type="checkbox"/>	☆	LinkShare	March Madness Sale! 50% Off All Sample Packs - Share Embe
<input type="checkbox"/>	☆	WESTERN UNION MONEY TR	Register Now: Social & Mobile Technologies Webinar - Social I
<input type="checkbox"/>	☆	Miss Beauty Musa	WESTERN UNION - Attn, We are grateful to contact you and anno
<input type="checkbox"/>	☆	American Musical Supply	Dearest - Dearest I know this mail will come to you as a surprise s
<input type="checkbox"/>	☆		Live Loud on Stage with Pro Gear up to 66% off - Speaker Syst

# 情感分析和观点挖掘

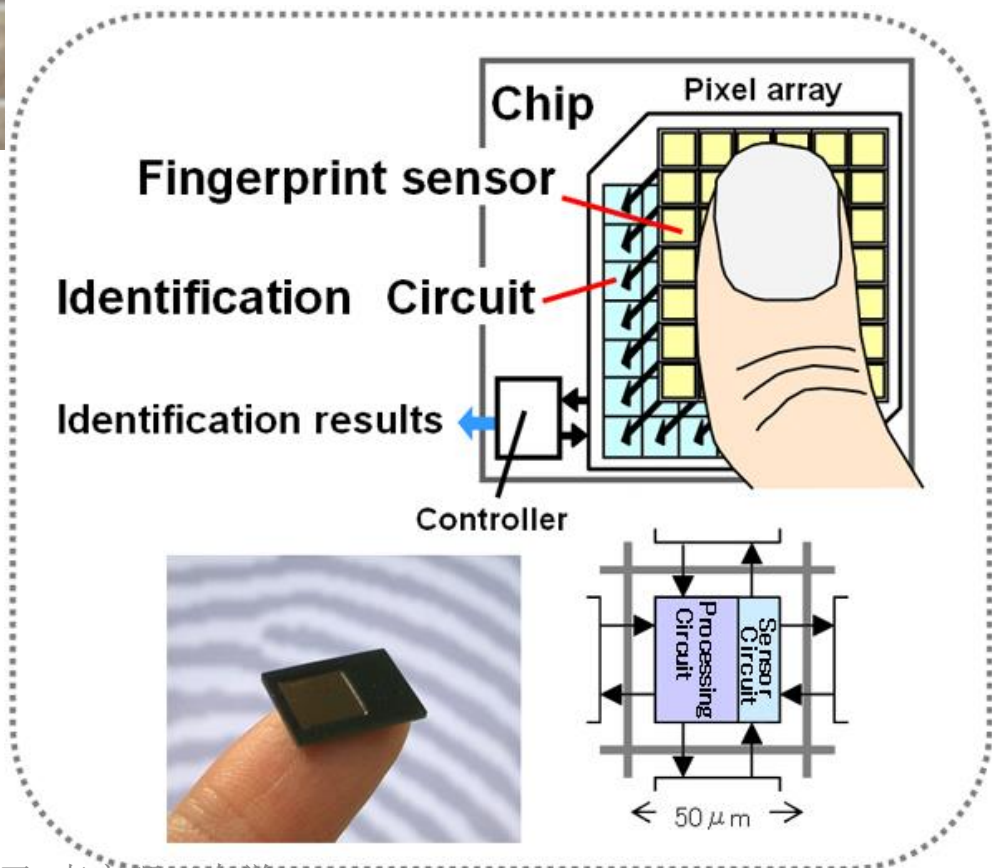
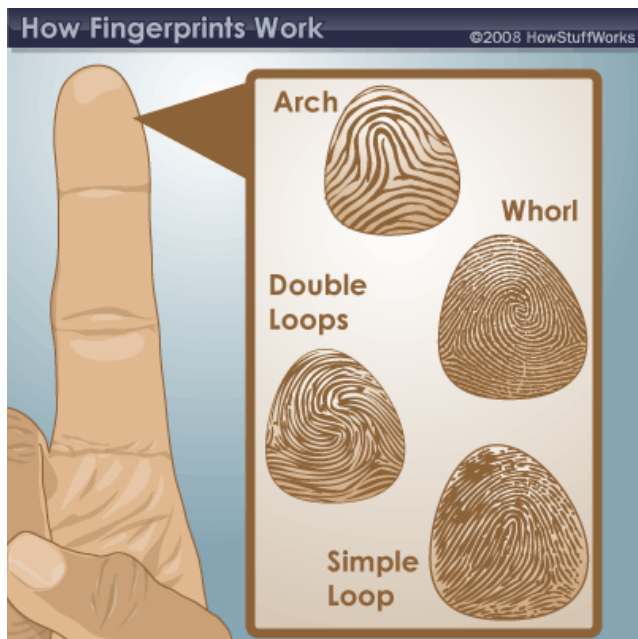
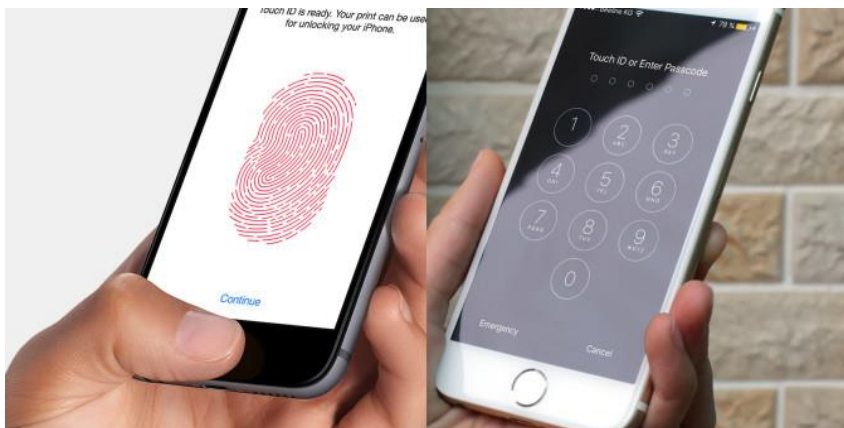


# 社交媒体分析





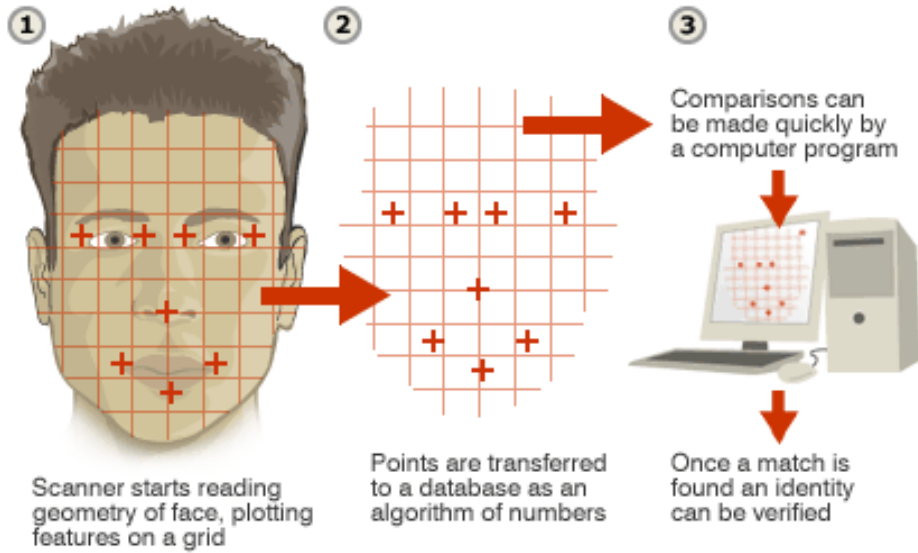
# 指纹识别



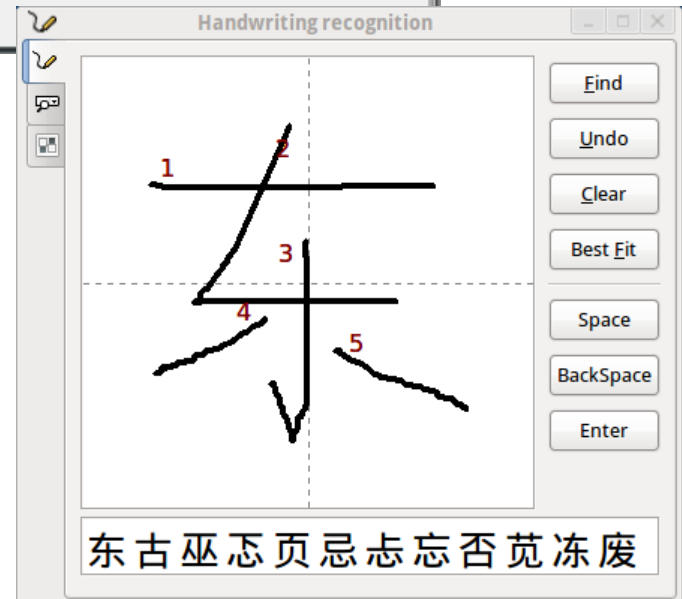
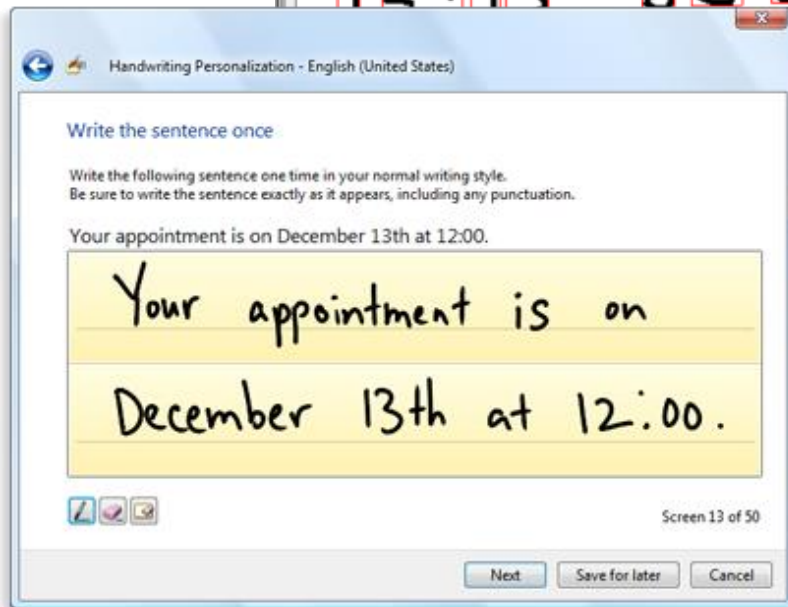
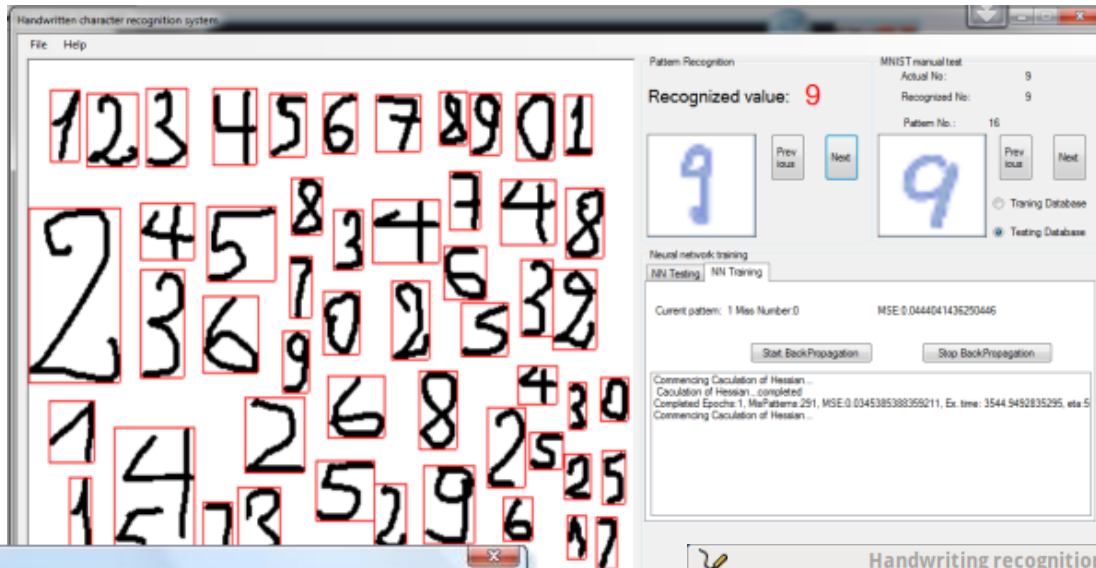
# 人脸识别



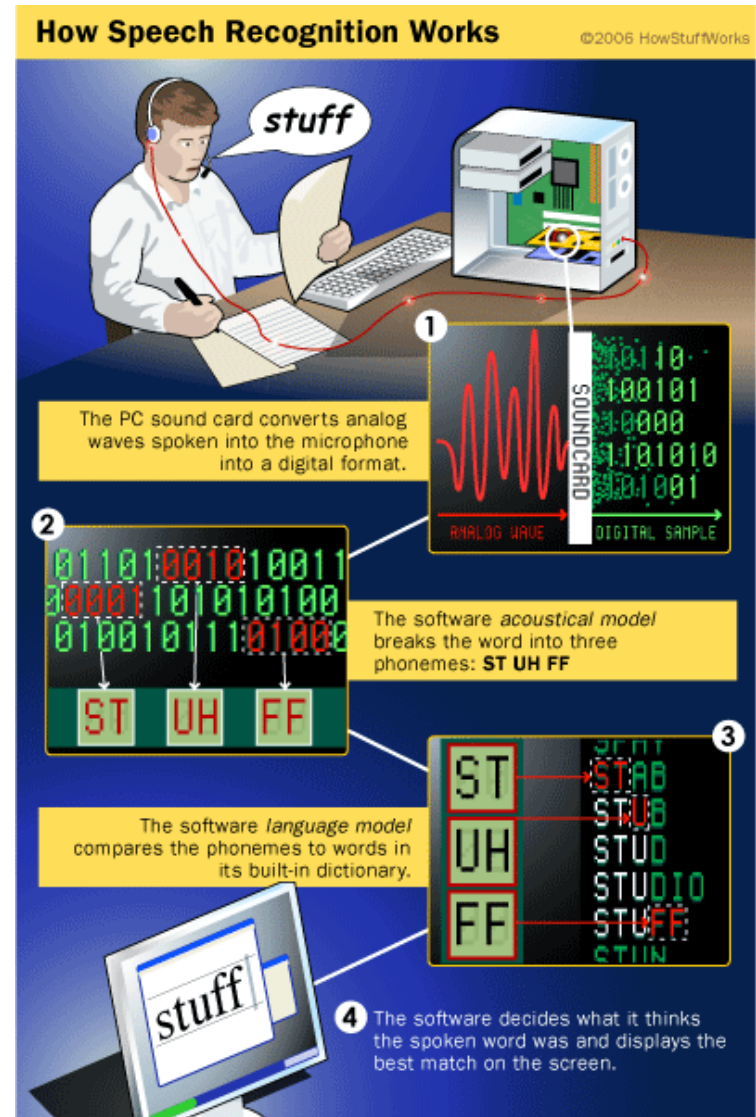
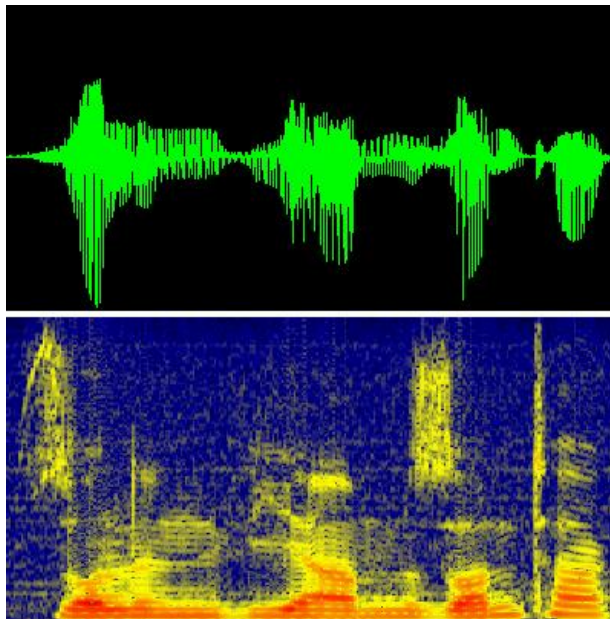
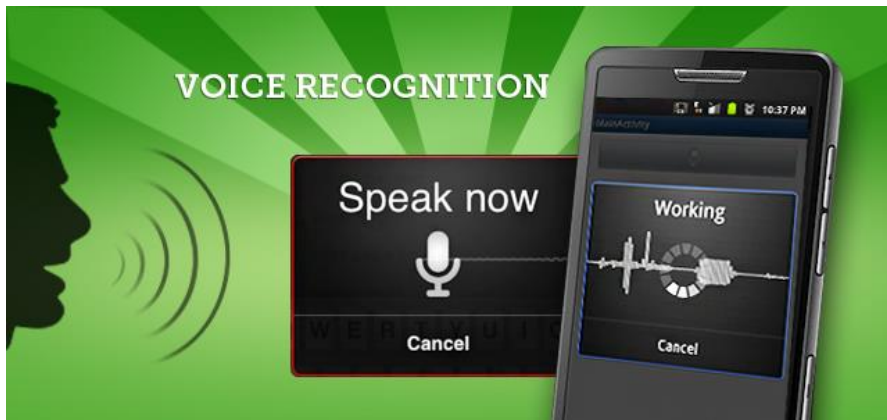
## HOW 2D FACIAL SCANNERS RECORD IDENTITIES



# 手写字符识别



# 语音识别



# 图片识别

中国移动 4G 18:36 39%

首页 识别结果



97% MATCH ▶



郁金香  
Didier's Tulip  
博爱体贴

中国移动 4G 18:41 36%

首页 识别结果



99% MATCH ▶



诸葛菜(二月兰)  
Orychophragmus Violaceus  
无私奉献

中国移动 4G 18:39 38%

首页 识别结果



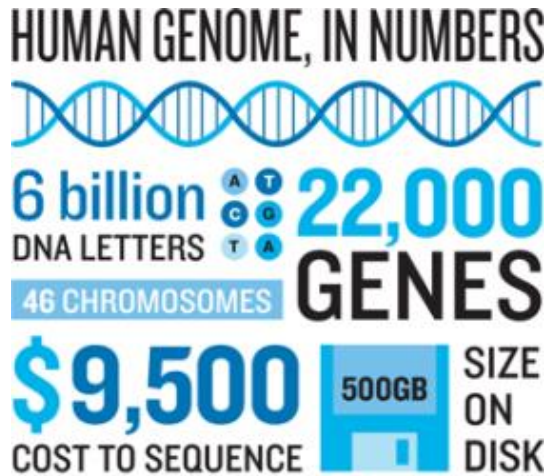
53% MATCH ▶



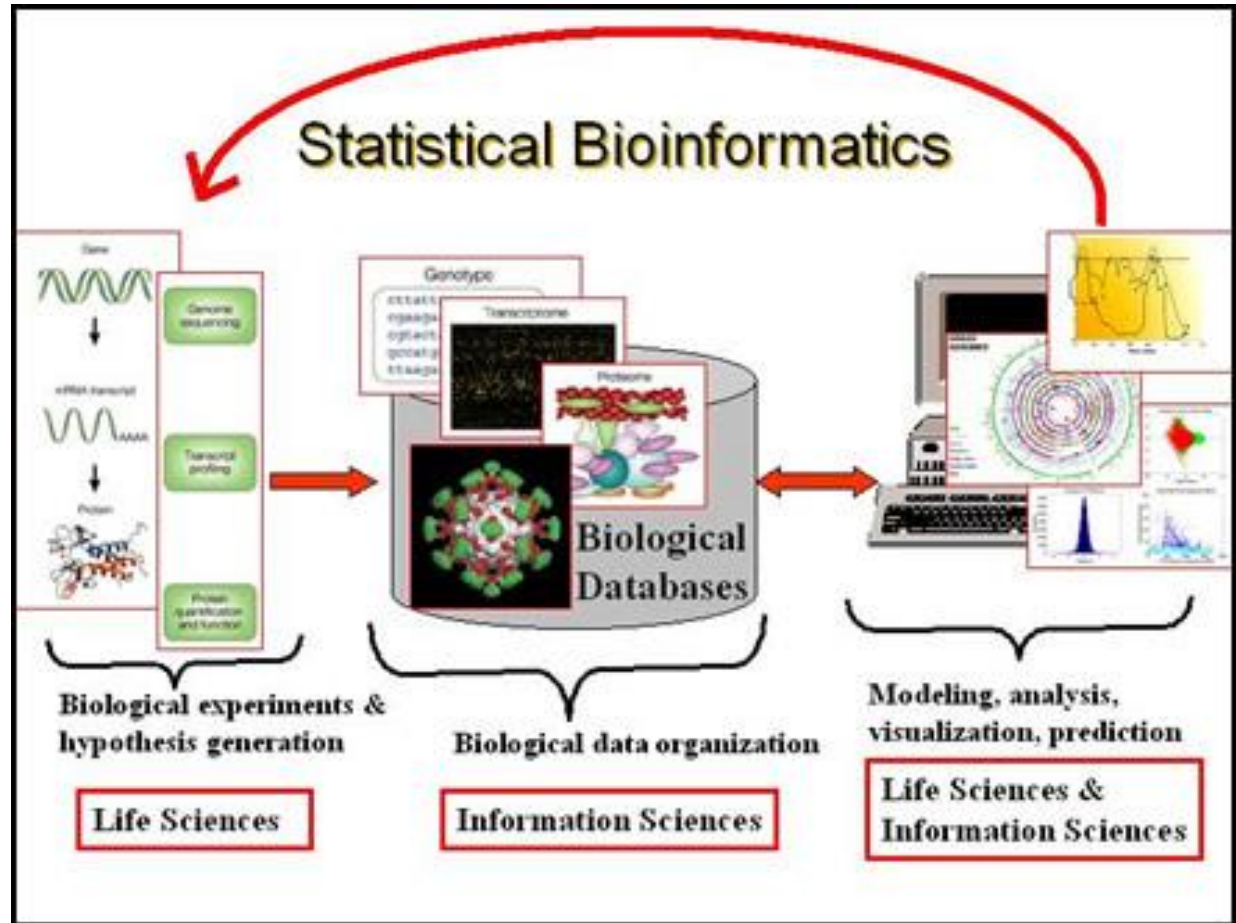
马蹄莲  
Calla Lily  
纯洁无瑕

<https://www.zhihu.com/question/51020471>

# 生物信息学



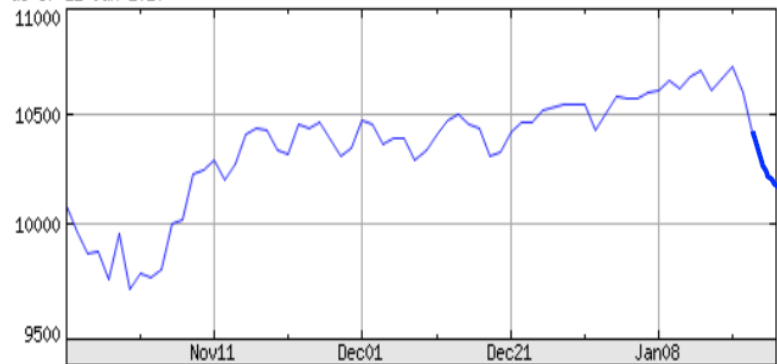
Sources: NIH, Illumina



# 股市预测



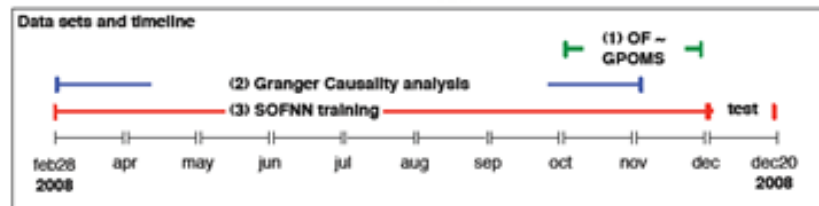
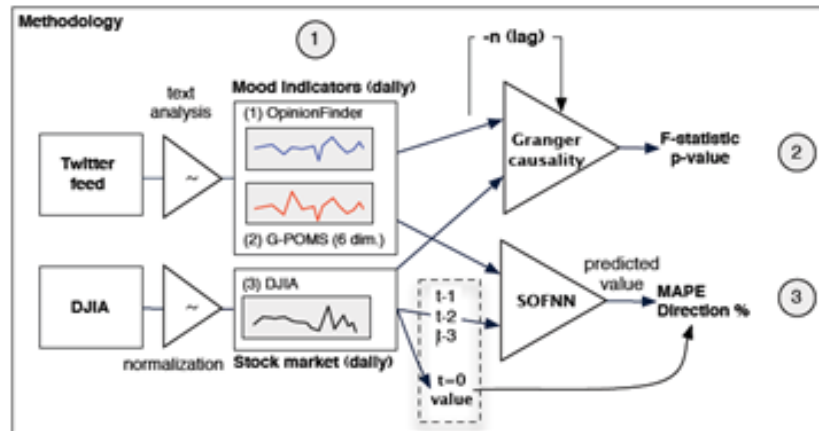
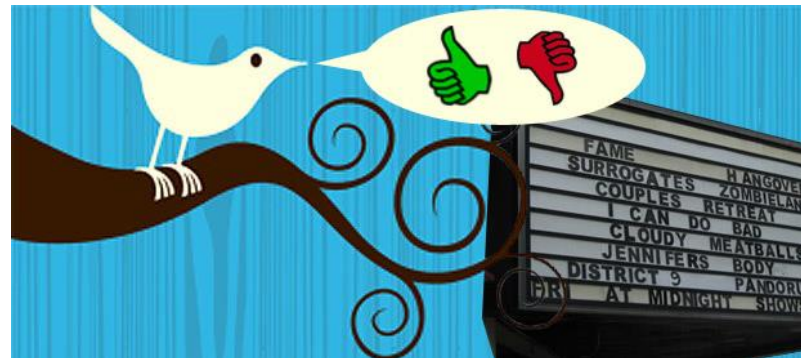
DJ INDU AVERAGE (DOW JONES & CO)  
as of 22-Jan-2010



Copyright 2010 Yahoo! Inc.

<http://finance.yahoo.com/>

X = Feb01



# 人类-机器 比赛



▶ 韓國棋王李世石三  
戰連敗 美聯社



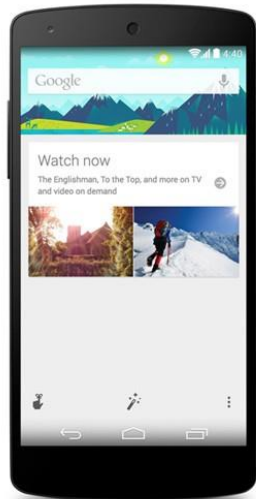


# 对话系统

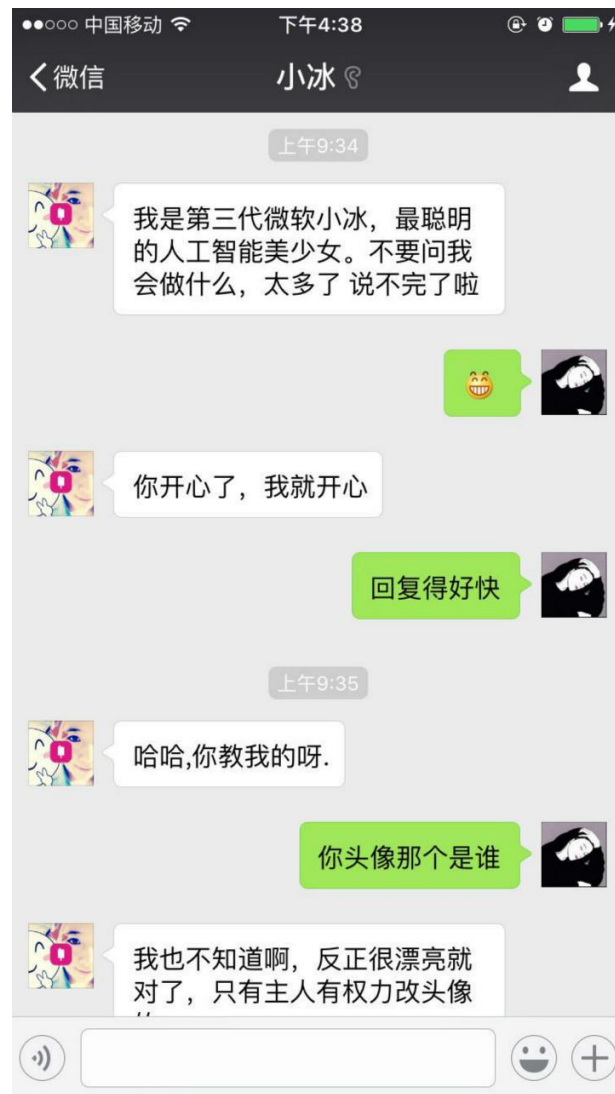
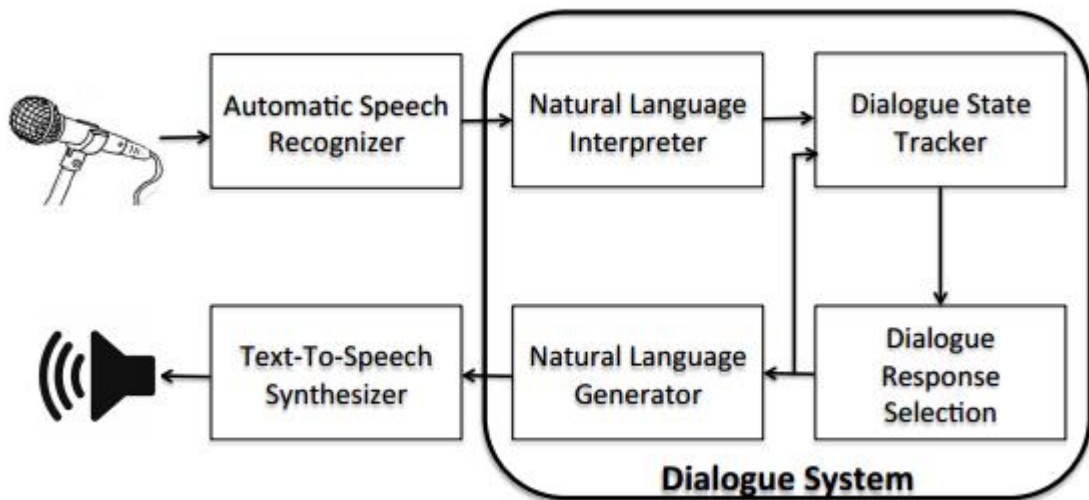
Apple Siri



Google Now



Windows Cortana



# 机器翻译



Translate

Turn off instant translation



Chinese English Spanish Detect language



English Chinese (Simplified) Spanish

Translate

## 【谷歌NMT，见证奇迹的时刻】

微信最近疯传人工智能新进展：谷歌翻译实现重大突破！值得关注和庆贺。mt 几乎无限量的自然带标数据在新技术下，似乎开始发力。报道说：

十年前，我们发布了 Google Translate（谷歌翻译），这项服务背后的核心算法是基于短语的机器翻译（PBMT:Phrase-Based Machine Translation）。

自那时起，机器智能的快速发展已经给我们的语音识别和图像识别能力带来了巨大的提升，但改进机器翻译仍然是一个高难度的目标。

今天，我们宣布发布谷歌神经机器翻译（GNMT Neural Machine Translation）系统，该系统使用先进的训练技术，能够实现到目前为止机器翻译的大提升。我们的全部研究成果详情请参阅我们《Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation》。

[Google NMT, witness the miracle of the moment]

Recent advances in microblogging crazy biography of artificial intelligence: Google translation to achieve a major breakthrough! Worthy of attention and celebration. Mt almost unlimited number of natural standard data in the new technology, it seems to start force. The report says:

Ten years ago, we released Google Translate, the core algorithm behind this service is PBMT: Phrase-Based Machine Translation.

Since then, the rapid development of machine intelligence has given us a great boost in speech recognition and image recognition, but improving machine translation is still a difficult task.

Today, we announced the release of the Google Neural Machine Translation (GNMT) system, which utilizes state-of-the-art training techniques to maximize the quality of machine translation so far. For a full review of our findings, please see our paper "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation."



# 无人车



# 更多应用

- 网络搜索和信息检索
- 自然语言理解
- 机器翻译
- 在图像和视频中定位/追踪/识别目标
- 金融预测和商业智能
- 医疗诊断，媒体图像分析
- 推荐系统
- ...

# 机器学习的关键词/概念

bayesian clustering conditional-distribution cost-function cross-entropy  
decision discriminative distribution em gaussian  
generative graphical-model inference joint-  
distribution least-square likelihood logistic-regression  
map ml model model-selection multinomial naive-bayes  
over-fitting predictive-function regression semi-supervised sequential-model  
supervised unsupervised

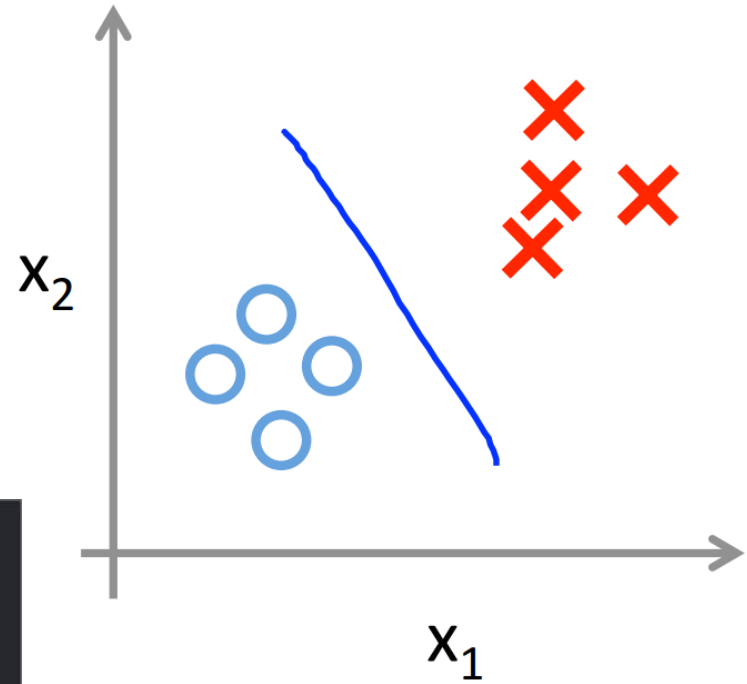
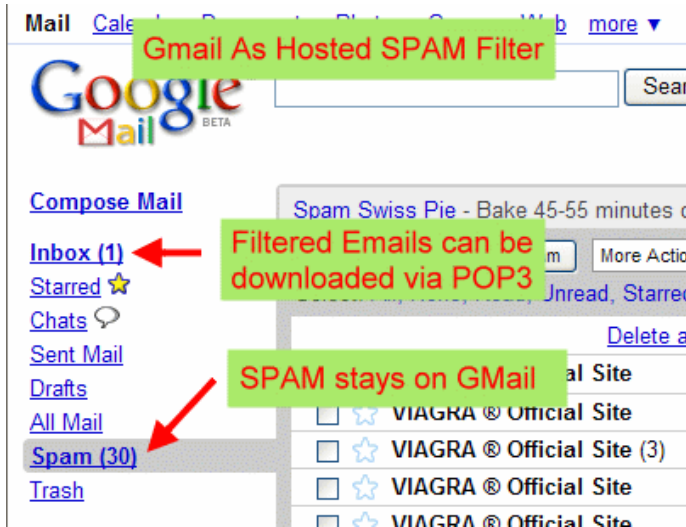
Tag-crowd of Bi-shop's PRML Book  
(模式识别与机器学习)

<http://research.microsoft.com/en-us/um/people/cmbishop/PRML>

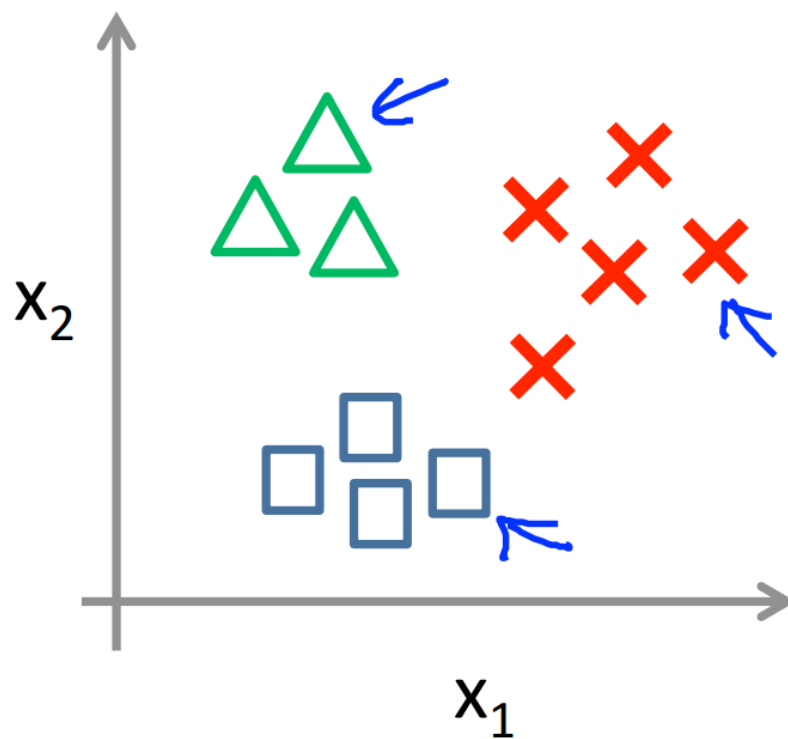
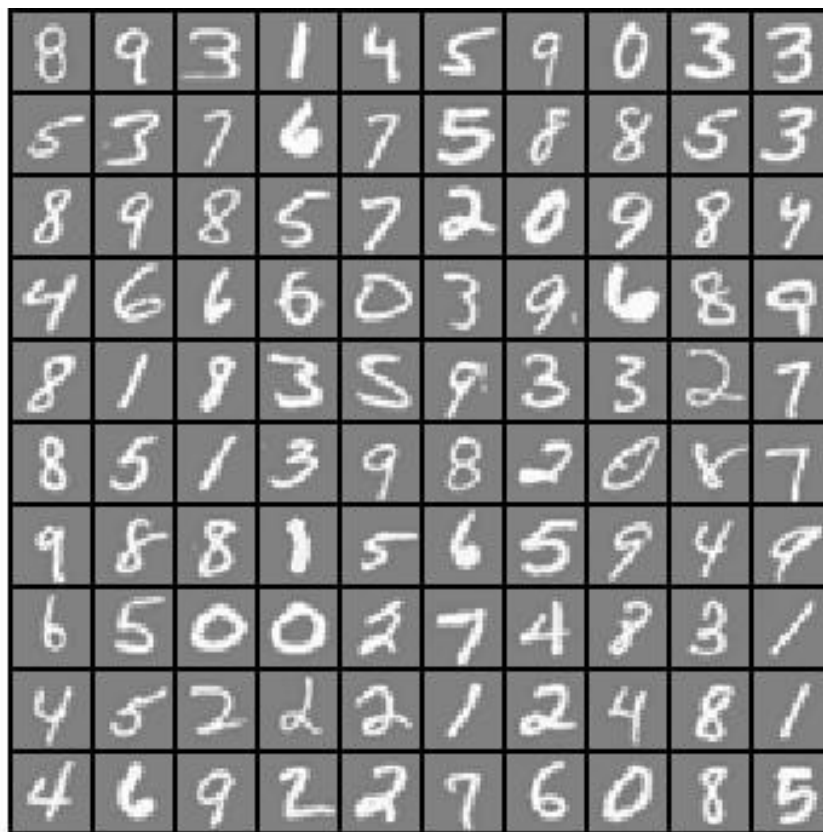
# 机器学习分类

- 监督学习: 给予一个输入对应一个输出结果的例子, 然后从新的输入预测输出结果。
  - 分类, 回归等
- 非监督学习: 给予一个输入对应一个输出结果的例子, 然后从新的输入预测输出结果
  - 聚类, 密度估计等.
- 半监督学习
- 集成学习
- 主动学习
- 迁移学习
- 增强学习
- 深度学习
- ...

# 二分类

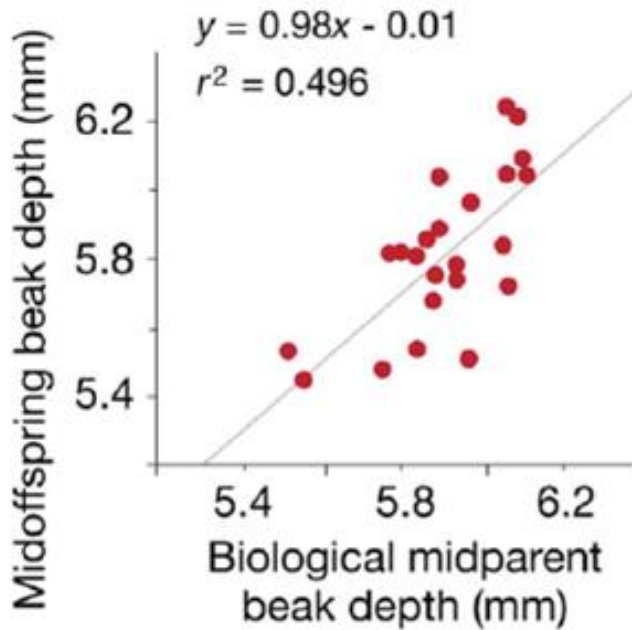


# 多分类





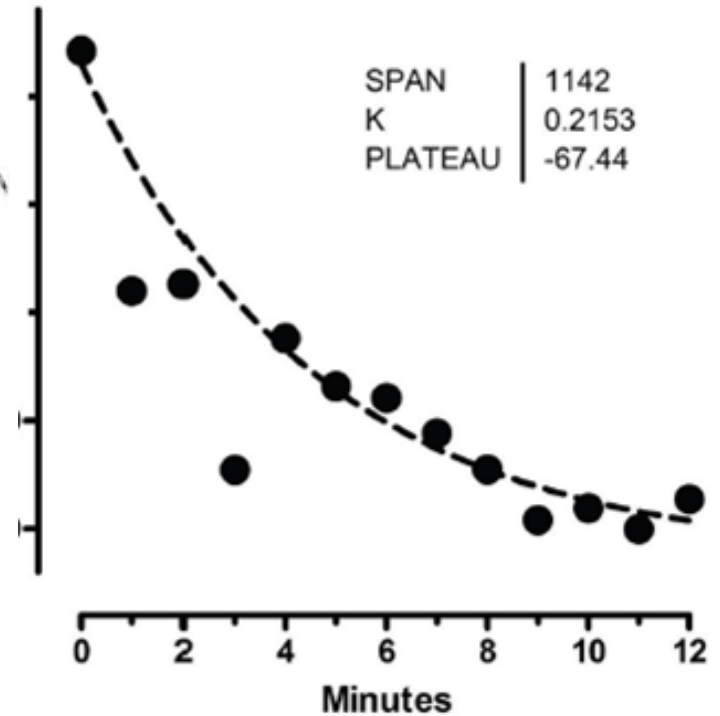
# 回归



线性回归



非线性回归

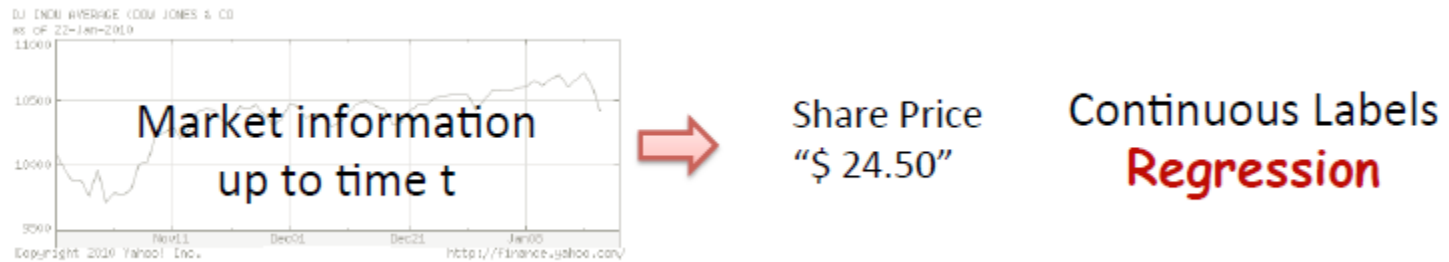


# 分类 vs. 回归

- 分类



- 回归



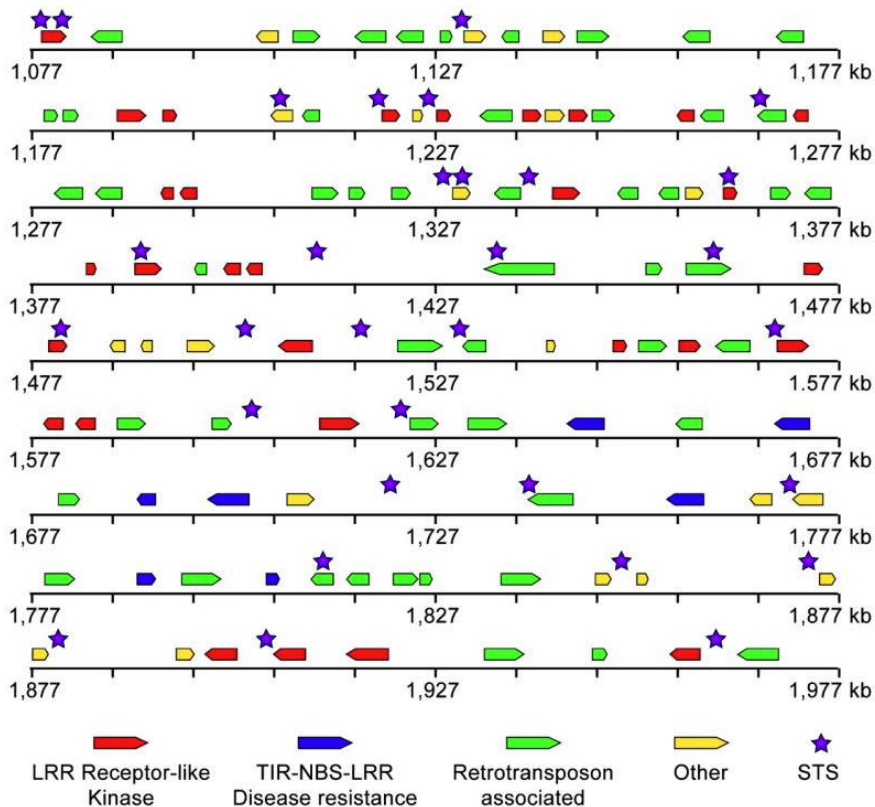
# 聚类

将相似的模型分组（单词，图像...）

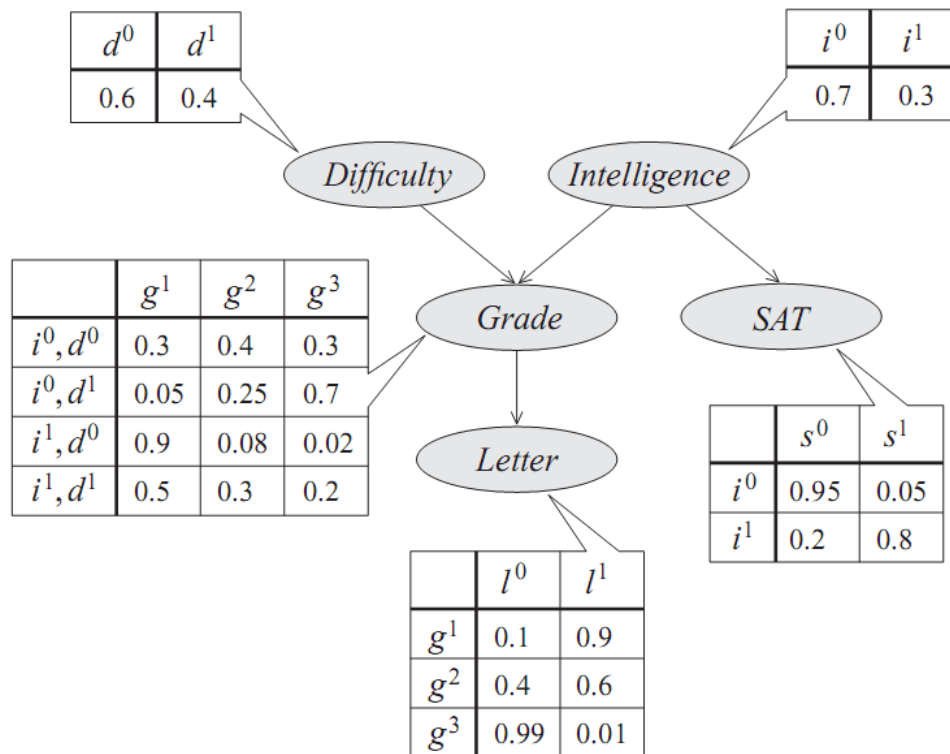


# 序列、图的学习

- Sequential 模型

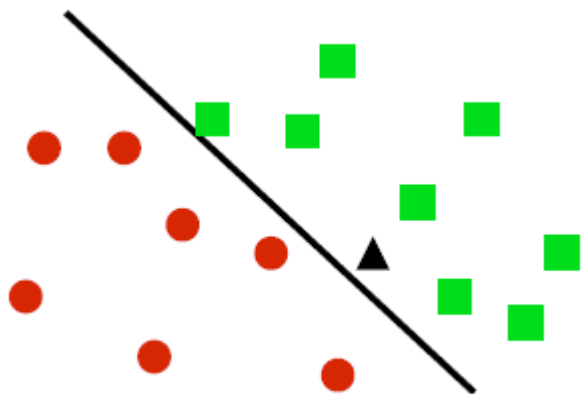


- Graphical 模型



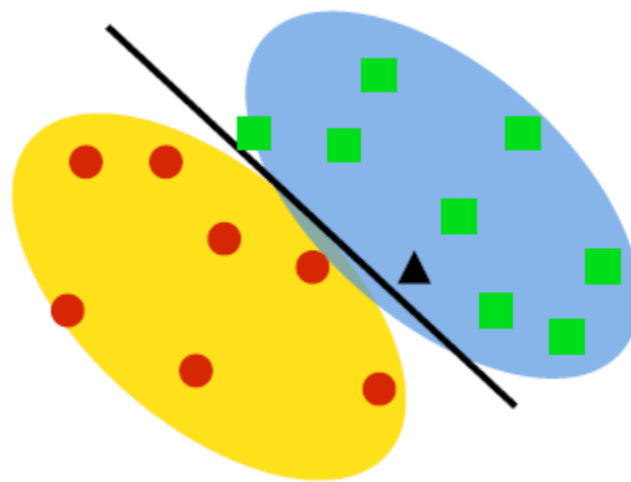
# 生成模型 vs. 判别模型

- 判别模型



用给定观察 $p(y|x)$ 标签的后验概率进行建模

- 生成模型



标签和观察  $p(x, y)$  的联合概率模型，然后用贝叶斯法则 $p(y|x) = p(x, y)/p(x)$ 来预测

# 机器学习阶段

- 训练阶段(用**训练数据**)

你可以从“黄金标准”中提供数据，并通过将输入与预期输出匹配来训练你的模型

- 测试阶段(用**测试数据**)

为了估计你的模型被训练的如何，并估计模型属性（如回归的平均误差，分类的准确性）

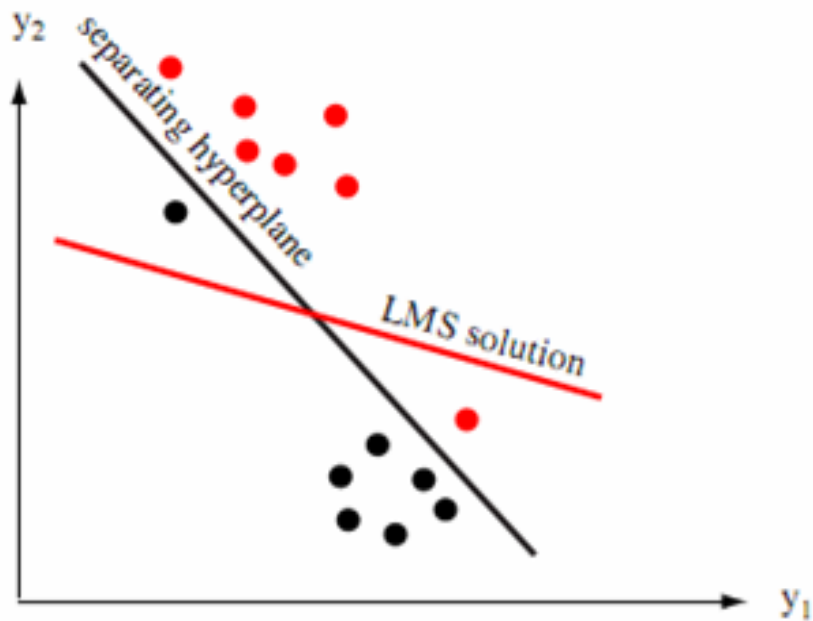
- 应用阶段(为**将来的数据**)

现在，你可以将新开发的模型应用于现实世界的的数据并获得结果

# 假设-学习-决策

- 假设
  - 具有（未知）参数（或结构）的数学模型
- 学习 (用来评估参数)
  - 最大似然估计（MLE）, MAP, 贝叶斯估计
  - 代价函数优化
- 决策
  - 贝叶斯决策规则
  - 直接预测函数

# 学习准则



- 感知机准则
- 最小均方(LMS)
- 最小交叉熵(CE)
- 最大间距准则
- 最大似然
- ...

哪一个线性超平面更好?

选择哪个学习准则?





# 优化方法

- 分析解决方案
- 梯度下降
- 随机梯度下降
- 牛顿法
- 拟牛顿法(BFGS)
- 有限内存BFGS (L-BFGS)
- 共轭梯度
- GIS
- IIS
- ...

# 基本的数学知识

- 微积分（微分，积分）
- 线性代数
- 概率论
- 优化方法



欢迎提问！