



# 数据贵过黄金?

Joseph Malkevitch / 文 曾铁勇 / 译

在不同的时期，不同的地方，黄金一直备受人类社会的推崇。2011年，某些地方黄金的价格每盎司突破了1900美元。多年来，黄金一直是对石油（目前每桶超过100美元）的重要性的形象比喻。但是数据有可能成为未来真正有价值的商品吗？如果是这样，这与数学有什么关系呢？

2012年的四月是数学宣传月，这一年的宣传主题是：数学、统计和海量数据。

本文的目的是要探讨短语“数据挖掘”的含义，以及为了促进这个领域的发展采用了哪些数学工具和思想。数据挖掘的起源之一是减小数据的占用空间，从而产生了对收集到的数据进行分析的需求。本文将给出一个非数学问题引起数学界关注后，相关的数学又如何发展的经典例子。人们不止一次看到，过去所发展的数学，只因为其“美”及智力上的吸引力，常常是洞悉新的应用环境的工具。在描述问题的背景后，我将给出一些引起广泛关注的和数据相关的例子，并探讨涉及数据挖掘的“人工智能”的方法。

## 无处不在的数据

什么是“数据”？作为多重领域使用的术语，数据最常代表的是“事实”或数值形式的统计。然而，有时这个术语又被用来表示将要被分析或用来做决断的信息，最近出现的对于数据的“定义”又涉及到用计算机处理或分析的数字、符号、字符串和表格。

虽然数据长伴我们，但利用数据进行更深入的探索并（或）影响决策的想法则相对较新。

几何和代数之根源可上溯千年，而“数据”数学则是近代产物，在过去的150年间发展迅猛。产生这种现象的部分原因是人们需要计算或分析大规模数据以获得重要信息，但这往往非常耗时，且受制于人为错误。因此，人们很自然地求助计算机来大规模采集和分析数据，以达到快速和准确的目的。

在看到不因噪音或偶然测量误差而出现的

模式这一意义下，要了解数据，需要懂得概率论。概率论和统计学是在许多方式下共存的两个科目。然而，概率是一个困难与微妙的学科。尽管概率论的数学基础相当坚实，但在数学之外的学科中用概率来数学建模，却异常复杂。当我们给出一个陈述：一枚硬币出现正面的机会稍大于出现背面的机会，比如假设出现正面的概率为 0.501，而出现背面的概率为 0.499，如何诠释这个陈述呢？如果在投掷方式不影响硬币出现正面或者背面的情况下（注意：某些技高一筹者可以多次投掷一枚硬币使得每次都出现正面）多次投掷硬币，将得到一个关于出现 H（代表正面）和 T（代表背面）的模式。例如，投掷一枚硬币 10 次而得到下面的模式：

TTHHHTTTTH

现在，基于这个有限的、很小的硬币投掷集合，出现正面的相对频率为 4/10（10 次投掷出现 4 次正面），而出现背面的相对频率为 6/10（10 次投掷出现 6 次背面）。由此你也许看到了这里出现的一个困难：对于任何一次固定数量的投掷，甚至是一次非常大数量的投掷，得到出现正面的相对频率为 0.501 和出现背面的相对频率为 0.499 的情况十分罕见！概率的“稳定的相对频率”的解释为：从长远来看，随着投掷次数变得越来越多，得到正面和背面的相对频率将分别为 0.501 和 0.499。不幸的是，似乎没有方法使得产生这种论点背后的直觉很“严谨”。除此之外，有些说法比如，这个电厂在未来 10 年将会发生核事故的概率为 0.00000001，或者明天在波士顿某些地区将会下雨的概率为 4/10 都没有很清晰的含义。

现代概率的观点是概率是受制于某些规则（公理）的系统，概率“直观的”性质产生于这个系统。这些规则意味着当有人使用“相对频率”的观点时，在一定意义上不会被误导。然而，这些年出现了其它的途径来“解释”概率的含义。显然，对于滚动骰子、投掷硬币及孩子的出生性别模式，可以较合理地理解概率意义下稳定的相对频率。然而，对于一个核电厂的燃料棒在未来 10 年内发生熔毁的可能性，概率的相对频率意味着什么呢？出现这类事件的历史非常短，所以对其稳定的相对频率的理解变得毫无意义。即使对于天气报告中经常出现的预报，比如明天下雨的可能性（概率）

是 80% 的预报，我们应该如何理解呢？

因数学家长期纠结于对概率这个概念所赋予的意义，导致了许多不同的“矛盾的”观点的出现。有鉴于此，概率和统计学的杰出贡献者伦纳德·萨维奇（Leonard Savage, 1917-1971）指出：“众所公认，统计学在某种程度上依赖于概率。但是，对于概率是什么以及它如何与统计学相联系，很少像今天如此激烈地争论并产生完全不同的观点。”



伦纳德·萨维奇（1917-1971）

问题之一是如何使用共同的语言（无论是英语、法语等）来表达不同环境下涉及到的“噪声”、“随机性”、“可能性”，或者是“意外”。放射性衰变机制显然不同于龙卷风会将袭击哪个州，或将会在哪天袭击这样的问题。

目前对于概率意义有多种诠释。作为一个例子，概率论的一种解释涉及到可根据经验获得的知识来帮助主观判断。这种方法也称为贝叶斯概率（即使是这个术语也存在几种不同的版本），它试图量化当和目标事件相关的某事件发生的概率正知后，如何来修正目标事件发生的概率。由这种角度产生的“概率”与数学家的正规方法产生的“概率”遵循相同的基本规则。然而，基于概率的不同表示所得到的推理方法也会有所不同。因此，会存在一些情况，因为采用不同的概率表示方法，做决断时就要从不同的可能中选择一种。估计这个复杂的议题将长期被数学家、统计学家和哲学家大范围地讨论。

## 概率和统计的先驱

在数学中，重要的思想凭空出世是非常罕见的。许多国家的科学家都为概率和统计的发展做出了贡献。本节将简要介绍一小部分重要贡献者。毫无疑问，概率论中相对频率的早期先驱之一是法国哲学家和数学家布莱士·帕斯卡（Blaise Pascal, 1623-1662）。帕斯卡的出发点是赌徒在实际赌局中的机会问题，在此前提下他提出了一些较深刻的见解。



布莱士·帕斯卡 (1623-1662)

英国牧师托马斯·贝叶斯（Thomas Bayes, 1702-1761）对概率论有了更深刻的理解，并且作出了极其重要的贡献。



托马斯·贝叶斯 (1702-1761)

贝叶斯的著名成果涉及到条件概率这个概念，这些条件概率可以针对于人们在实验中看到的或者在某些假设下产生的。这样实验或假设出现的结果称为事件。当无偏地投掷硬币（这时出现正面或背面的相对频率假定为  $1/2$ ）10 次，那么 10 次中正好出现 7 次背面的概率是多少，或者 10 次中正好出现 7 次正面的概率是多少？假设生男生女的概率为  $1/2$ ，以及不同的出生事件是彼此独立的（即一个小孩的出生不影响其他小孩的出生），那么，一对夫妇生的前两个小孩都是女孩的概率是多少？如果有顺序的字符串  $GGB$  表示第一个小孩是女孩，第二个小孩是女孩，第三个小孩是男孩，那么所询问的概率可表示为  $P(GG)$  是多少？我们也可以利用角标来表示出生的顺序，如  $B_1, G_2, G_1B_2$  分别代表三个事件：第一个小孩是男孩，第二个是女孩，第一个是女孩且第二个是男孩。现在假设我们知道第一个小孩是女孩。问两个小孩都是女孩的概率是多少？这里我们问的是在给定第一个小孩是女孩的条件下，两个小孩都是女孩的“条件概率”。更一般地，令  $P(X|Y)$  = 事件  $Y$  发生的条件下事件  $X$  发生的概率，也可以记  $P(Y|X)$  = 事件  $X$  发生的条件下事件  $Y$  发生的概率。不难看出，直观上假设  $P(Y)$  不为零，则

$$P(X|Y) = \frac{P(X \text{ 且 } Y)}{P(Y)}$$

利用两个集合  $X$  和  $Y$  的交集的符号，上式可以转变为

$$P(X|Y) = \frac{P(X \cap Y)}{P(Y)}$$

经过片刻的思考，我们就知道当  $P(X)$  和  $P(Y)$  都不为零时， $P(X|Y)$  和  $P(Y|X)$  并不需相等。

例如，对于上面的小孩出生问题，我们可以知道  $P(G_1G_2|G_1) = 1/2$ ,  $P(G_1|G_1G_2) = 1$ ；这里我们用到的下面的事实： $P(G_1G_2) = 1/4$ ,  $P(G_1) = 1/2$ 。

在计算条件概率时，贝叶斯给出了著名的贝叶斯定理，这个定理在贝叶斯死后才得以

发表。直观上理解，贝叶斯定理提供了一个框架，可以从中排序出对导致事件发生的“因素”（事件）有“相对”影响力的事件。更具体地说，假设我们有一系列事件  $X_1, X_2, \dots, X_k$ ，事件间相互排斥，即只可以有一个事件发生。因此，在同样一个空间中，一个实验的结果属于由这些事件组成的集合。假定  $E$  是一个概率不为零的事件。假设我们知道  $E$  一定发生，那么如何计算  $P(X_i|E)$ ？贝叶斯定理是一个“公式”，由这个公式可以计算出  $P(X_i|E)$  的值，即已知  $E$  发生，一个特定因素  $X_i$  发生的概率。

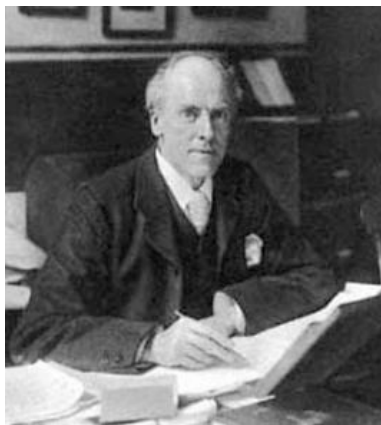
数学史的较早期，即 18 世纪或者更早时，几乎所有对数学做出重大贡献的人同时也是伟大的物理学家。如牛顿 (Isaac Newton, 1643-1727)，欧拉 (Leonhard Euler, 1707-1783)，拉普拉斯 (Pierre-Simon Laplace, 1749-1827)，勒让德 (Adrien-Marie Legendre, 1752-1833)，高斯 (Karl Friedrich Gauss, 1777-1855)，他们不仅是杰出的数学家同时也是杰出的物理学家。作为物理学家，他们对数据中的“噪声”，即物理量的测量误差，均有所关注。统计和概率论的某些理论正是源于这种考虑。上面提到的数学家高斯以不同方式思考这个问题，由此产生了“最小二乘法”。正态分布的相关思想亦由此产生。另外，统计学不仅关注自然科学，而且关注社会科学。在这方面的先驱者之一是比利时的数学家和科学家阿道夫·凯特勒 (Adolphe Quetelet, 1796-1874)，他利用正态曲线来研究人类的特点，并从分析的角度来研究犯罪规律。其工作属于将统计思想用于社会学的早期实践。



阿道夫·凯特勒 (1796-1874)

在现代，越来越多的人对统计及统计与概率论的结合做出贡献。在这里，我将要展示的不只是来自不同背景不同国家的人们对统计和概率的贡献，更要说明有多少成果根源于近代。当然了，把这个话题说透彻大概需要一本书的内容。

卡尔·皮尔逊 (Karl Pearson, 1857-1936) 在英格兰出生及去世。皮尔逊曾就读于剑桥大学，而他的大部分职业生涯在伦敦大学的大学院度过。皮尔逊对利用数学的思想来研究进化感兴趣，这也促使他提出了新的统计思想和方法。他 1894 年提出的标准偏差这个术语为众多文科修读统计的学生所知，标准偏差也成为度量数据发散程度的通用术语。皮尔逊还促进了使用大样本的统计检验思想的发展。



卡尔·皮尔逊 (1857-1936)

像所有的数学分支一样，在精细的审视下，统计学的发展有着丰富而复杂的历史，其推动者的身份不仅仅是数学家，他们同时也从事其他学术领域的智力活动。凯恩斯 (John Maynard Keynes, 1883-1946) 在经济学领域大名鼎鼎，然而他在剑桥大学学习了数学，并于 1921 年出版了关于概率论的一本重要专著。凯恩斯的工作被数学家、哲学家罗素勋爵关注，同时也引起了弗兰克·拉姆齐 (Frank Ramsey) 的注意。除了组合数学方面的著名成就 (今天我们称为拉姆齐定理)，拉姆齐同样在概率论和统计学方面做出了重要的贡献。最初凯恩斯的概率观点趋向于稳定的相对频率方法，但随着时间的推移，他的观点逐渐趋向于主观性更强的信念体系，可能在经济学中这个更有吸引力吧。



约翰·梅纳德·凯恩斯 (1883-1946)

罗纳德·费舍尔 (Ronald Fisher, 1890-1962) 是另一个曾受教于剑桥大学的英国数学家，他在统计学研究中硕果累累。所谓的 F 检验即以他的名字命名。费舍尔同时也对实验设计感兴趣。他在英国罗森斯德 (Rothamsted) 的一个农业试验站工作多年，在那里他参与了一个实验，研究不同的培育方式对植物生长的影响。培育方式包括浇水制度、土壤的不同类型、施肥方案等。利用统计分析这一工具，可以找出不同的培育方式对不同类型的植物所产生的影响，并可根据结果改进方案来增加粮食的产量。很遗憾的是，费舍尔和皮尔逊在他们职业生涯中曾一度因费舍尔的统计思想和方法发生过激烈的争执。



罗纳德·费舍尔 (1890-1962)

耶日·奈曼 (Jerzy Neyman, 1894-1981) 出生于俄罗斯，但后期研究主要在美国进行。期间他也曾在伦敦住过一段时间，在那里，他与卡尔·皮尔逊的儿子、统计学家埃贡·皮尔逊 (Egon Pearson 1895-1980) 相互交流。后来，奈曼在加州大学伯克利分校工作，并使该校的统计系闻名世界。



耶日·奈曼 (1894-1981)

布鲁诺·德费奈蒂 (Bruno de Finetti, 1906-1985) 出生于奥地利，在意大利接受教育，最后在罗马去世。他以提倡从“主观”的角度理解可能性的含义而闻名。



布鲁诺·德费奈蒂 (1906-1985)

随着计算机功能不断增强，数学家开始借助这个强大工具来探索和发现数据中隐藏的信息。约翰·图基 (John Tukey, 1915-2000) 是这一领域的先驱之一，他创造了二进制数字中的术语“比特”。图基一开始是化学系的学生，

但最终从普林斯顿获得数学博士学位。他长期在贝尔实验室工作，并曾担任实验室副主任。为了从数据中获得尽可能多的信息，图基和他的同事研究如何开发计算机的新功能、增大内存及加快计算机的处理速度。人眼对视觉模式异常敏感，所以图基和贝尔实验室的同事还探讨了如何利用人类视觉系统来显示数据以探索数据。他最传世的工作应该是他和合作者詹姆斯·库利（James Cooley）于1965年提出的快速傅里叶变换，这是信号处理领域最重要的工具之一。



约翰·图基（1915-2000）

## 数据密集型课题

有时我们会称最近在数据里“溺水”了。出现这个说法的部分原因是非常多的数据正在被生成、收集及存储，可我们大多数人并没有时间去浏览这些数据，更不要说去思考数据里蕴含的意义。从某种意义上来说，21世纪美国人生活的各个方面都在被数据驱动。为了更清楚理解这一点，下面将列举几个“数据密集”的领域。

即使是在高速计算机出现之前，某些领域已经在广泛地收集数据了。下面会简短给出几个这样的领域，科学家不断研究新的统计方法，以便对付这些领域中的大数据问题。

### 天气

获得好的天气预报通常要求解数学中的偏微分方程，而这些方程需要数据作为输入。不同来源的详细信息收集到一起，可以从中获得精确的天气信息。人们对天气一直比较感兴趣，这导致了大量数据的生成。这不仅是因为天气关系到上班时带不带伞的问题，而且有助于庄稼种植和庄稼收割等重要民生问题。与天气同时收集的信息有气压、温度、相对湿度、风力强度和方向、雨雪量等。这些数据都可用来做决策或规划，比如利用这些数据可以分析冰川和极地冰盖的融化及这些变化对临海国家的影响。

### 医学成像

医学成像数量及使用量急剧增多，许多成像的突破都受益于数学的支持。例如，约翰·拉东

（Johann Radon）的工作，即拉东变换，使得层析成像成为可能。其他技术，包括计算机速度和存储，也在促进医学成像的发展。随着医学信息系统的不断完善，在治疗同一个病人时，医生们可以共享诸如CT和MRI扫描及验血所获得的数据。这些数据可以帮助医生正确诊断症状复杂的病人。

### 新药物的开发

开发安全且有效的新药物非常重要，这时数学工具也常被用来理解疾病治疗的本质，特别是遗传病和传染病的治疗。例如，数学已经被用来为不同阶段的艾滋病人寻找合适的治疗方案。人们知道，随着时间的推移，一些药物（抗生素和抗疟疾药物）的作用会逐渐减小，这是因为引起疾病的病原体随着药物的使用产生抗体。所以人们需要不断寻找新的方法来调整现有的药物，以提高药效，或在病原体产生抗体之前更新药物。从病人处收集到的数据及药物的结构形式，是用数学和统计思想来研制更好的药物治疗方案的原始材料。

### 基因组学 / 计算生物学 / 生物信息学

克里克（Francis Crick）和沃森（James Watson）提出了著名的遗传学模型，为生物学和数学同时开辟了一个崭新的领域，这不仅吸引了统计与概率，亦吸引了其他数学分支。概括地说，作为间接遗传分子，DNA可以看成是由四个字母ACGT组成的序列，这里每个字母都是一个

特定的核苷酸的简称。每个个体都有自己特殊的由这些字母序列组成的基因。基于 DNA 模型，每个物种都有一个区别于其他物种的特定的基因组成，不同的物种具有不同数量的染色体及染色体基因组的数量。越来越快的序列分析仪，产生了越来越多的关于个体基因组和不同物种基因组的数据。生物学家声称地球上有一百万个物种，但有人说目前地球上的物种的数量远大于这个数字。关于物种的正确定义的部分问题是统计的问题。判别两个对象是否属于同一类，可设计一些基于对象间一系列的“度量”作为距离。给定一些 DNA 序列，我们可以利用各种“距离”来度量这些序列是否相似。在此环境下发展的统计研究，诞生了生物信息学这个全新的领域。

### 宇宙学 / 天文学

人们在远离灯光的晴朗的夜晚用肉眼就可以看到数量令人惊叹的星星。现在，使用安装在山顶的望远镜（越来越多的安装在南半球，因为那里人造光源比较少）、人造卫星、太空望远镜（特别是哈勃望远镜），海量的关于天体的数据正在被收集。利用成像技术和统计思想，天文学家正在尝试更深入地理解我们宇宙的本质，并寻找适合人类生存的星球。这与物理学家试图了解时间、空间和引力的本质工作相辅相成。

### 赌博与投资

在赌博或投资追逐利益的过程中，大量的数据生成。纽约证券交易所（NYSE）每天都产生大量的数据。在越来越短的时间跨度内，成百上千种股票被追踪并试图预测其规律，以便吸引更多股民在这些交易所投资。许多与股票打交道的人被称为技术分析师。这些人研究股票价格的变化趋势（通常是用复杂的统计工具），再结合其他的一些信息以期投资获利。金融数学中的新思想与统计学这个重要工具，已经被用来辅助一些国家在影响最低的控制商业周期，并用来帮助他们更好地理解世界金融市场。

### 粒子物理

在美国和欧洲，用来加速原子运动速度的机器已建成多年，包括著名的费米实验室的粒子加速器和欧洲核子研究组织（CERN）的大型强子对撞机。建造这种机器的目的是为了帮助

物理学家了解物质的本质以及组成物质的粒子。虽然所谓的标准模型已经可以成功地解释大部分关于物质的原理，但仍有许多物理学家关注的谜团无法解开。谜团之一是希格斯玻色子是否存在。虽然可以找到这种粒子的能量水平的范围在减小，但是在各种粒子加速器的大量实验中还未发现这种粒子（2012年7月，欧洲核子研究组织称其发现了与希格斯玻色子特征基本相符的新粒子。2013年3月14日，欧洲核子研究组织宣布，先前探测到的新粒子是希格斯玻色子。——编者注。）每次用这些加速器做实验时，通常会有大量的图像和数据生成。数据挖掘技术被用来辅助检查这些数据，来寻找代表意想不到的事物的个别“粒子衰变”，或寻找可以支持我们对已知现象的理解的证据。

### 交通信息

住在城市的居民非常关心能否准时上下班或能否准时参加约会。这时，没有比因交通堵塞而迟到更令人沮丧的事情了。目前正在开发的许多系统通过配置传感器和相机来检测交通流量，并试图利用由这些系统产生的数据给司机实时的建议，让他们合理地规划路线以便减少塞车的时间。这些数据同样有助于了解在一座桥上如何设置收费亭，以及决定在有双向车道的桥上如何合理分配车道。

### 教育

作为监控教育行业发展的手段之一，数据收集愈发重要。学校，不论是中小学、大学、研究生院还是专业学院，都要平衡学生的数量和学生的质量。各级政府都努力使公立学校有足够的承载力，并通过数据判断这些学校是否满足社会的需求。

### 电子邮件和互联网

现在越来越多的人花大量的时间查阅电子邮件或上网。从个人或商业角度看，查电邮或上网都产生了大量数据。目前许多公司都在关注收集用户或群体（通常是基于用户的邮政编码来得到）的信息，用来改进浏览器的搜索结果。有些公司甚至会花钱使自己出现在网页搜索中显眼的位置上，而数据也会告诉商家广告某些特定的搜索词可能会给公司带来收益。社会学家、经济学家也对电子邮件和互联网相关数据收集到的信息感兴趣。

## 数据挖掘的工具

数据挖掘的基础是统计数学，同时又受益于人工智能和计算学习理论（computational learning theory）。

人工智能领域已经吸引了数学家、计算机学家和哲学家的关注。严格地说，这个领域是指由人类设计和编程的计算机在“创造”和“思考”的时候可以表现出类似人类的能力。

历史上，人工智能的目标是令计算机可以玩人类需要投入高层次思考的游戏，这些游戏强调技巧而不存在运气的因素，如跳棋、象棋和围棋。例如，在人工智能的早期历史上，有人相信当计算机“武装”上优秀玩家用到的规则就可以下象棋，然而这种方法并不十分成功。随着计算机拥有越来越快的处理速度和越来越大的内存，人工智能开始借助计算机的这些特征再往前发展。今天，对于象棋棋盘上的任何一个给定的位置，计算机可以找到所有可能的走法以及对手相应的走法，利用这个“暴力”方法，结合“位置评估”可找到最佳的移动和应对方法，正因如此，IBM的深蓝计算机终于击败了一位世界冠军。

最近，IBM设计了一个名为沃森（Watson）的计算机系统，它击败了著名智力抢答节目Jeopardy中最成功的人类选手，这个节目使参与者在语言诙谐的气氛下挑战一系列复杂而广泛的知识。选手选择一个具有特定分数的主题，问题的难度与所附分数的值成正比。参与者不时地选择问题，如果答对了，则使收到的金钱加倍。有时对于一个特殊类别中的问题，选手可以拿已得到的金钱中的一部分来对正确答案“下注”。因此，这里有一个“终极危险”的回合，如果选手确信自己给出的是正确答案，则可以尝试通过拿已得的部分资金下注以期超过

对手。2011年2月14至16日，在黄金时间播出的一系列特别电视节目里，尽管在回答问题时出现一些奇怪的“行为”，IBM的沃森计算机“系统”打败了两名非常优秀的人类对手。沃森的信号手的技能（当它准备好回答问题时就按铃）非常好，但是当主持人叫它回答问题必须简短的时候其反应就会差一些。

计算学习理论主要是指探索可以使机器针对一个具体的目标改进自身性能算法。例如，处理大量电子邮件的人的一个困扰是垃圾邮件的增长。垃圾邮件包含通过电子邮件发送的无用信息，没有人愿意收到它们。已经开发出软件用于将有可能的垃圾邮件和推销信息放到一个目录下，如果在一定的时间段内没有被查看，则这些邮件会被自动删除。用户通常可以设置垃圾邮件过滤器的参数，以便尽可能地避免删除那些用户愿意阅读的邮件。利用贝叶斯统计思想，可以设计出越来越成功的垃圾邮件过滤器，即使在邮箱的用户没有指出过滤器哪部分的决策不满意的条件下也很成功。

理解数据的一种方法是找到可以解释这些数据的数学模型。这方面一个成功的例子是利用均值和标准偏差来对数据的集中趋势和变化程度建模。实际上，无论是从物理实验、心理学实验收集到的数据，还是英语教授观察莎士比亚戏剧中句子的长短，都可以用同样的数学工具来分析。另外一个展现数学魅力的例子是概率密度函数和概率分布函数的思想的产生，这使得探讨数据是服从正态分布还是服从其它分布成为可能。然而，生成数据的数量越来越大，以及在哪些领域数据可以发挥更大的作用，对我们依然是重大的挑战。数学家、统计学家和计算机科学家正在努力应对这一挑战。

原文链接：<http://www.ams.org/samplings/feature-column/fc-2012-04>

注：译者为曾铁勇博士，香港浸会大学