

我经常在 TopLanguage 讨论组上推荐一些书籍，也经常问里面的牛人们搜罗一些有关的资料，人工智能、机器学习、自然语言处理、知识发现（特别地，数据挖掘）、信息检索 这些无疑是 CS 领域最好玩的分支了（也是互相紧密联系的），这里将最近有关机器学习和人工智能相关的一些学习资源归一个类：

首先是两个非常棒的 Wikipedia 条目，我也算是 wikipedia 的重度用户了，学习一门东西的时候常常发现是始于 wikipedia 中间经过若干次 google，然后止于某一本或几本著作。

第一个是“[人工智能的历史](#)”（History of Artificial Intelligence），我在讨论组上写道：

而今天看到的这篇文章是我在 wikipedia 浏览至今觉得最好的。文章名为《人工智能的历史》，顺着 AI 发展时间线娓娓道来，中间穿插无数牛人故事，且一波三折大气磅礴，可谓“事实比想象更令人惊讶”。人工智能始于哲学思辨，中间经历了一个没有心理学（尤其是认知神经科学的）的帮助的阶段，仅通过牛人对人类思维的外在表现的归纳、内省，以及数学工具进行探索，其间最令人激动的是 Herbert Simon（决策理论之父，诺奖，跨领域牛人）写的一个自动证明机，证明了罗素的数学原理中的二十几个定理，其中有一个定理比原书中的还要优雅，Simon 的程序用的是启发式搜索，因为公理系统中的证明可以简化为从条件到结论的树状搜索（但由于组合爆炸，所以必须使用启发式剪枝）。后来 Simon 又写了 GPS（General Problem Solver），据说能解决一些能良好形式化的问题，如汉诺塔。但说到底 Simon 的研究毕竟只触及了人类思维的一个很小很小的方面——Formal Logic，甚至更狭义一点 Deductive Reasoning（即不包含 Inductive Reasoning, Transductive Reasoning (俗称 analogic thinking)）。还有诸多比如 Common Sense、Vision、尤其是最为复杂的 Language、Consciousness 都还谜团未解。还有一个比较有趣的就是有人认为 AI 问题必须要以一个物理的 Body 为支撑，一个能够感受这个世界的物理规则的身体本身就是一个强大的信息来源，基于这个信息来源，人类能够自身与时俱进地总结所谓的 Common-Sense Knowledge（这个就是所谓的 Embodied Mind 理论。），否则像一些老兄直接手动构建 Common-Sense Knowledge Base，就很傻很天真了，须知人根据感知系统从自然界获取知识是一个动态的自动更新的系统，而手动构建常识库则无异于古老的 Expert System 的做法。当然，以上只总结了很小一部分我个人觉得比较有趣或新颖的，每个人看到的有趣的地方不一样，比如里面相当详细地介绍了神经网络理论的兴衰。所以我强烈建议你看自己一遍，别忘了里面链接到其他地方的链接。

顺便一说，[徐宥](#)同学打算找时间把这个条目翻译出来，这是一个相当长的条目，看不动 E 文的等着看翻译吧:)

第二个则是“[人工智能](#)”（Artificial Intelligence）。当然，还有[机器学习](#)等等。从这些条目出发能够找到许多非常有用和靠谱的深入参考资料。

然后是一些书籍

书籍:

1. 《**Programming Collective Intelligence**》，近年出的入门好书，培养兴趣是最重要的一环，一上来看大部头很容易被吓走的:P
2. Peter Norvig 的《**AI, Modern Approach 2nd**》（无争议的领域经典）。
3. 《**The Elements of Statistical Learning**》，数学性比较强，可以做参考了。
4. 《**Foundations of Statistical Natural Language Processing**》，自然语言处理领域公认经典。
5. 《**Data Mining, Concepts and Techniques**》，华裔科学家写的书，相当深入浅出。
6. 《**Managing Gigabytes**》，信息检索好书。
7. 《**Information Theory: Inference and Learning Algorithms**》，参考书吧，比较深。

相关数学基础（参考书，不适合拿来通读）:

1. 线性代数：这个参考书就不列了，很多。
2. 矩阵数学：《**矩阵分析**》，Roger Horn。矩阵分析领域无争议的经典。
3. 概率论与统计：《概率论及其应用》，威廉·费勒。也是极牛的书，可数学味道太重，不适合做机器学习的。于是讨论组里的 **Du Lei** 同学推荐了《**All Of Statistics**》并说到

机器学习这个方向，统计学也一样非常重要。推荐 All of statistics，这是 CMU 的一本很简洁的教科书，注重概念，简化计算，简化与 Machine Learning 无关的概念和统计内容，可以说是很好的快速入门材料。

4. 最优化方法：《**Nonlinear Programming, 2nd**》非线性规划的参考书。《**Convex Optimization**》凸优化的参考书。此外还有一些书可以参考 wikipedia 上的最优化方法条目。要深入理解机器学习方法的技术细节很多时候（如 SVM）需要最优化方法作为铺垫。

王宁同学推荐了好几本书:

《**Machine Learning, Tom Michell**》, 1997.

老书，牛人。现在看来内容并不算深，很多章节有点到为止的感觉，但是很适合新手（当然，

不能"新"到连算法和概率都不知道)入门。比如决策树部分就很精彩,并且这几年没有特别大的进展,所以并不过时。另外,这本书算是对97年前数十年机器学习工作的大综述,参考文献列表极有价值。国内有翻译和影印版,不知道绝版否。

《**Modern Information Retrieval, Ricardo Baeza-Yates et al**》. 1999

老书,牛人。貌似第一本完整讲述 IR 的书。可惜 IR 这些年进展迅猛,这本书略有些过时了。翻翻做参考还是不错的。另外, Ricardo 同学现在是 Yahoo Research for Europe and Latin America 的头头。

《**Pattern Classification (2ed)**》, Richard O. Duda, Peter E. Hart, David G. Stork

大约也是01年左右的大块头,有影印版,彩色。没读完,但如果想深入学习 ML 和 IR,前三章(介绍,贝叶斯学习,线性分类器)必修。

还有些经典与我只有一面之缘,没有资格评价。另外还有两本小册子,论文集性质的,倒是讲到了不少前沿和细节,诸如索引如何压缩之类。可惜忘了名字,又被我压在箱底,下次搬家前怕是难见天日了。

(呵呵,想起来一本:《**Mining the Web – Discovering Knowledge from Hypertext Data**》)

说一本名气很大的书:《**Data Mining: Practical Machine Learning Tools and Techniques**》。Weka 的作者写的。可惜内容一般。理论部分太单薄,而实践部分也很脱离实际。DM 的入门书已经不少,这一本应该可以不看了。如果要学习了解 Weka,看文档就好。第二版已经出了,没读过,不清楚。

信息检索方面, **Du Lei** 同学再次推荐:

信息检索方面的书现在建议看 Stanford 的那本《**Introduction to Information Retrieval**》,这本书刚刚正式出版,内容当然 up to date。另外信息检索第一大牛 Croft 老爷也正在写教科书,应该很快就要面世了。据说是非常 practical 的一本书。

对信息检索有兴趣的同学,强烈推荐翟成祥博士在北大的暑期学校课程,这里有全 slides 和阅读材料: <http://net.pku.edu.cn/~course/cs410/schedule.html>

maximzhao 同学推荐了一本机器学习:

加一本书: Bishop, 《**Pattern Recognition and Machine Learning**》. 没有影印的,但是网上能下到。经典中的经典。Pattern Classification 和这本书是两本必读之书。《Pattern Recognition and Machine Learning》是很新(07年),深入浅出,手不释卷。

最后，关于人工智能方面（特别地，决策与判断），再推荐两本有意思的书，

一本是《**Simple Heuristics that Makes Us Smart**》

另一本是《**Bounded Rationality: The Adaptive Toolbox**》

不同于计算机学界所采用的统计机器学习方法，这两本书更多地着眼于人类实际上所采用的认知方式，以下是我在讨论组上写的简介：

这两本都是德国 ABC 研究小组（一个由计算机科学家、认知科学家、神经科学家、经济学家、数学家、统计学家等组成的跨学科研究团体）集体写的，都是引起领域内广泛关注的书，尤其是前一本，后一本则是对 Herbert Simon（决策科学之父，诺奖获得者）提出的人类理性模型的扩充研究），可以说是把什么是真正的人类智能这个问题提上了台面。核心思想是，我们的大脑根本不能做大量的统计计算，使用 fancy 的数学手法去解释和预测这个世界，而是通过简单而鲁棒的启发法来面对不确定的世界（比如第一本书中提到的两个后来非常著名的启发法：再认启发法（*cognition heuristics*）和选择最佳（*Take the Best*）。当然，这两本书并没有排斥统计方法就是了，数据量大的时候统计优势就出来了，而数据量小的时候统计方法就变得非常糟糕；人类简单的启发法则充分利用生态环境中的规律性（*regularities*），都做到计算复杂性小且鲁棒。

关于第二本书的简介：

1. 谁是 [Herbert Simon](#)

2. 什么是 [Bounded Rationality](#)

3. 这本书讲啥的：

我一直觉得人类的决策与判断是一个非常迷人的问题。这本书简单地说可以看作是《决策与判断》的更全面更理论版本。系统且理论化地介绍人类决策与判断过程中的各种启发式方法（*heuristics*）及其利弊（为什么他们是最优化方法在信息不足情况下的快捷且鲁棒的逼近，以及为什么在一些情况下会带来糟糕的后果等，比如学过机器学习的都知道朴素贝叶斯方法在许多情况下往往并不比贝叶斯网络效果差，而且还速度快；比如多项式插值的维数越高越容易 *overfit*，而基于低阶多项式的分段样条插值却被证明是一个非常鲁棒的方案）。

在此提一个书中提到的例子，非常有意思：两个团队被派去设计一个能够在场上接住抛过来的棒球的机器人。第一组做了详细的数学分析，建立了一个相当复杂的抛物线近似模型（因为还要考虑空气阻力之类的原因，所以并非严格抛物线），用于计算球的落点，以便正确地接到球。显然这个方案耗资巨大，而且实际运算也需要时间，大家都知道生物的神经网络中生物电流传输只有百米每秒之内，所以 *computational complexity* 对于生物来说是个宝贵资源，所以这个方案虽然可行，但不够好。第二组则采访了真正的运动员，听取他们总结自己

到底是如何接球的感受，然后他们做了这样一个机器人：这个机器人在球抛出的一开始一半路程啥也不做，等到比较近了才开始跑动，并在跑动中一直保持眼睛于球之间的视角不变，后者就保证了机器人的跑动路线一定会和球的轨迹有交点；整个过程中这个机器人只做非常粗糙的轨迹估算。体会一下你接球的时候是不是眼睛一直都盯着球，然后根据视线角度来调整跑动方向？实际上人类就是这么干的，这就是 **heuristics** 的力量。

相对于偏向于心理学以及科普的《决策与判断》来说，这本书的理论性更强，引用文献也很多而经典，而且与人工智能和机器学习都有交叉，里面也有不少数学内容，全书由十几个章节构成，每个章节都是由不同的作者写的，类似于 **paper** 一样的，很严谨，也没啥废话，跟《**Psychology of Problem Solving**》类似。比较适合 **geeks** 阅读哈。

另外，对理论的技术细节看不下去的也建议看看《决策与判断》这类书（以及像《别做正常的傻瓜》这样的傻瓜科普读本），对自己在生活中做决策有莫大的好处。人类决策与判断中使用了很多的 **heuristics**，很不幸的是，其中许多都是在适应几十万年前的社会环境中建立起来的，并不适合于现代社会，所以了解这些思维中的缺点、盲点，对自己成为一个良好的决策者有很大的好处，而且这本身也是一个非常有趣的领域。

（完）

P.S. 大家有什么好的资料请至[讨论组上留言](#)。