

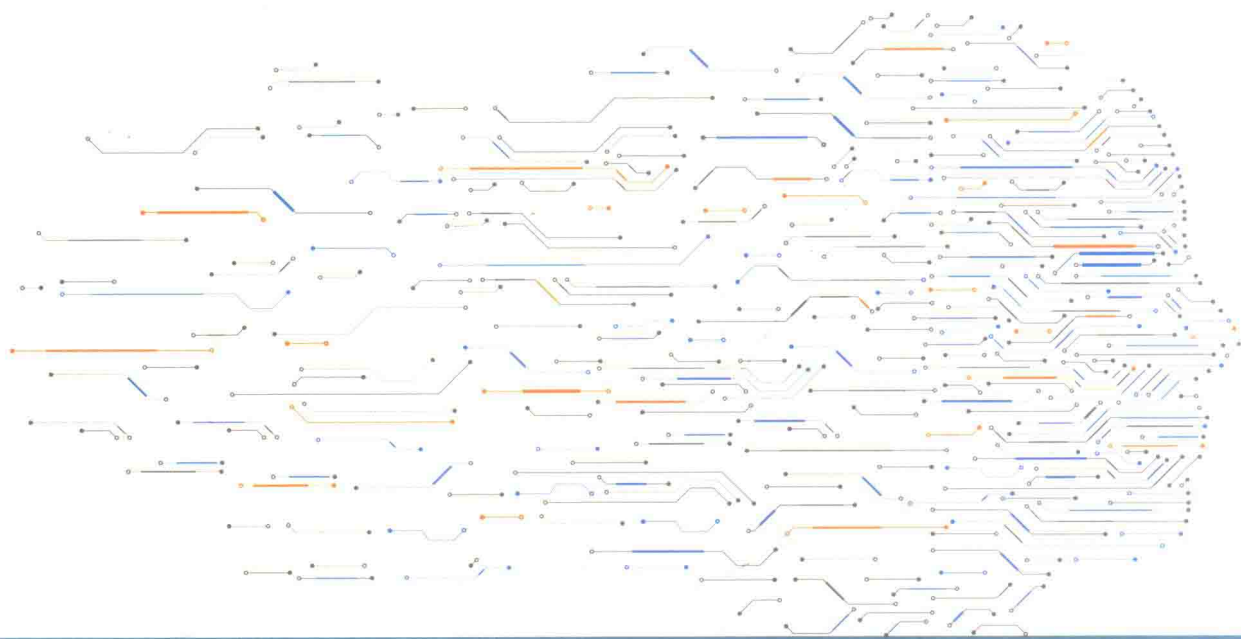
Distributed Machine Learning

Theories, Algorithms, and Systems

分布式机器学习

算法、理论与实践

刘铁岩 陈薇 王太峰 高飞 © 著



机械工业出版社
China Machine Press

智能科学与技术丛书

Distributed Machine Learning

Theories, Algorithms, and Systems

分布式机器学习

算法、理论与实践

刘铁岩 陈薇 王太峰 高飞 著



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

分布式机器学习: 算法、理论与实践 / 刘铁岩等著. —北京: 机械工业出版社, 2018.9
(2018.11 重印)
(智能科学与技术丛书)

ISBN 978-7-111-60918-6

I. 分… II. 刘… III. 机器学习 IV. TP181

中国版本图书馆 CIP 数据核字 (2018) 第 220552 号

分布式机器学习: 算法、理论与实践

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码: 100037)

责任编辑: 姚 蕾 迟振春

责任校对: 张惠兰

印 刷: 中国电影出版社印刷厂

版 次: 2018 年 11 月第 1 版第 2 次印刷

开 本: 185mm × 260mm 1/16

印 张: 17.25

书 号: ISBN 978-7-111-60918-6

定 价: 89.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88379426 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzit@hzbook.com

版权所有 · 侵权必究

封底无防伪标均为盗版

本书法律顾问: 北京大成律师事务所 韩光 / 邹晓东

最近几年，机器学习在许多领域都获得了前所未有的成功，由此也彻底改变了人工智能的发展方向，引发了大数据时代的到来。其中最富有挑战性的问题是由分布式机器学习解决的。所以，要了解机器学习究竟能够带来什么样前所未有的新机遇、新突破，就必须了解分布式机器学习。

相比较而言，机器学习这个领域本身是比较单纯的领域，其模型和算法问题基本上都可以被看成纯粹的应用数学问题，而分布式机器学习则不然，它更像一个系统工程，涉及数据、模型、算法、通信、硬件等许多方面，这更增加了系统了解这个领域的难度。刘铁岩博士和他的合作者的这本书，从理论、算法和实践等多个方面对这个新的重要学科给出了系统、深刻的讨论。这无疑是雪中送炭，这样的书籍在现有文献中还难以找到。对我个人而言，这也是我早就关注但一直缺乏系统了解的领域，所以看了这本书，我也是受益匪浅。相信对众多关注机器学习的工作人员和学生，这也是一本难得的好书。

我是2012年在我组织的“数据科学与信息产业”会议上认识铁岩的。后来虽然见面不多，但我一直关注他的工作。他和合作者在百忙之中抽出宝贵的时间来写这本书，对整个机器学习、大数据和人工智能领域都是很大的贡献。相信他们的辛勤劳动会得到行业的回报。

鄂维南

2018年6月

如果说人工智能技术将造就人类的未来时代，那么作为人工智能的核心支撑，机器学习将会像电力一样无处不在。事实上，机器学习现在已经炙手可热，不仅学界关注、业界聚焦、政府重视，甚至在街头巷尾也常有所闻。回望十几年前很多人还以为机器学习是机械类专业内容，恍如隔世。

机器学习备受关注的的原因之一，是它已经在众多现实应用中发挥了巨大作用，尤其在若干困难任务上带来了超出一般预料的成功。于是，人们热情高涨，对于以机器学习为核心的智能产业的前景无限憧憬，而如何让机器学习技术在业界的大规模任务中更充分地发挥威力，则成为热议的话题。

业界的大规模机器学习任务往往涉及如何充分地利用“大数据”、如何有效地训练“大模型”。使用价格昂贵的高性能设备，例如 TB 级内存的计算服务器未尝不可，但硬件能力的增长速度显然比不上机器学习所面对数据的增长速度，因此目前业界更主流的解决方案是分布式机器学习。

分布式机器学习并非分布式处理技术与机器学习的简单结合。一方面，它必须考虑机器学习模型构成与算法流程本身的特点，否则分布式处理的结果可能失之毫厘、谬以千里；另一方面，机器学习内含的算法随机性、参数冗余性等，又会带来一般分布式处理过程所不具备的、宜于专门利用的便利。

值得一提的是，市面上关于机器学习的书籍已有许多，但是分布式机器学习的专门书籍还颇少见。

刘铁岩博士是机器学习与信息检索领域的国际著名专家，带领的微软亚洲研究院机器学习研究团队成果斐然。此次他们基于分布式机器学习方面的丰富经验推出《分布式机器学习：算法、理论与实践》一书，将是希望学习和了解分布式机器学习的中文读者的福音，必将有力促进相关技术在中国的推广和发展。

周志华

于南京

2018年6月

近年来，人工智能取得了飞速的发展，实现了一个又一个技术突破。这些成功的幕后英雄是海量的训练数据、超大规模的机器学习模型以及分布式的训练系统。一系列有关分布式机器学习的研究工作，从并行模式、跨机通信到聚合机制，从算法设计、理论推导到系统构建，都在如火如荼地展开。人们不仅发表了大量的学术论文，也开发出一批实用性很强的分布式机器学习系统。本书的目的是向读者全面展示分布式机器学习的现状，深入分析其中的核心技术问题，并且讨论该领域未来发展的方向。本书既可以作为研究生从事分布式机器学习方向研究的参考文献，也可以作为人工智能从业者进行算法选择和系统设计的工具书。

全书共 12 章。第 1 章是绪论，向大家展示分布式机器学习这个领域的全景。第 2 章介绍机器学习的基础知识，其中涉及的基本概念、模型和理论，会为读者在后续章节中更好地理解分布式机器学习的各项技术奠定基础。第 3 章到第 8 章是本书的核心部分，向大家细致地讲解分布式机器学习的框架及其各个功能模块。其中第 3 章对整个分布式机器学习框架做综述，而第 4 章到第 8 章则针对其中的数据与模型划分模块、单机优化模块、通信模块、数据与模型聚合模块分别加以介绍，展示每个模块的不同选项并讨论其长处与短板。接下来的三章是对前面内容的总结与升华。其中第 9 章介绍由分布式机器学习框架中不同选项所组合出来的各式各样的分布式机器学习算法，第 10 章讨论这些算法的理论性质（例如收敛性），第 11 章则介绍几个主流的分布式机器学习系统（包括 Spark MLlib、Multiverso 参数服务器系统和 TensorFlow 数据流系统）。最后的第 12 章是全书的结语，在对全书内容进行简要总结之后，着重讨论分布式机器学习这个领域未来的发展方向。

有关本书的写作，因为涉及分布式机学习的不同侧面，不同的章节对读者预备知识的要求有所不同。尤其是涉及优化算法和学习理论的部分，要求读者对于最优化理论和概率统计有一定的知识储备。不过，如果读者的目的只是

熟悉主流的分布式机器学习框架和系统，则可以跳过这些相对艰深的章节，因为其余章节自成体系，对于理论部分没有过多的依赖。

我仍然清晰地记得，两年以前华章公司的姚蕾编辑多次找到我，希望我能撰写一本关于分布式机器学习的图书。一方面被姚蕾的诚意打动，另一方面也考虑到这样一本书对于在校研究生和人工智能从业者可能有所帮助，我最终欣然应允。然而，平时工作过于繁忙，真正可以用来写书的时间非常有限，所以一晃就是两年的时光，直至今日本书才与读者见面，内心十分惭愧。

回顾这两年的写作过程，有很多人需要感谢。首先，我要感谢本书的联合作者：陈薇博士负责书中与优化算法和学习理论有关的内容，王太峰和高飞则主要负责通信机制、聚合模式和分布式机器学习系统等方面的内容。没有他们夜以继日的努力，本书无法成文。在写作过程中，本书的各位作者得到了家人的大力支持。写书之路实属不易，如果没有她（他）们的默默奉献，作者们很难集中精力，攻克这个艰巨的任务。其次，我要感谢诸多为本书的写作做出过重要贡献的人：我在中国科学技术大学的博士生郑书新花费了大量的精力和时间帮助我们整理了全书的参考文献；北京大学的孟琪同学则帮助我们在全书做了细致的校验；华章公司的编辑姚蕾和迟振春对我们的书稿提出了很多宝贵的意见；普林斯顿大学教授、中国科学院院士鄂维南博士，以及南京大学教授周志华博士分别为本书题写了推荐序。正是因为这么多幕后英雄的奉献，本书才得以顺利面世。最后，我还要感谢微软亚洲研究院院长洪小文博士，他的大力支持使得我们在分布式机器学习这个领域做出了很多高质量的研究工作，也使得我们有机会把这些成果记录下来，编纂成书，与更多的同行分享。

惭愧的是，即便耗时两载，即便集合了多人的智慧和努力，本书的写作仍然略显仓促。加之分布式机器学习这个领域飞速发展，本书成稿之时，又有很多新的研究成果发表，难以周全覆盖。再则，本书的作者才疏学浅，书中难免有疏漏、错误之处，还望读者海涵，不吝告知，日后加以勘误，不胜感激。

刘铁岩

于北京中关村

2018年6月

刘铁岩 微软亚洲研究院副院长。刘博士的先锋性研究促进了机器学习与信息检索之间的融合，被国际学术界公认为“排序学习”领域的代表人物。近年来在深度学习、分布式学习、强化学习等方面也颇有建树，发表论文 200 余篇，被引用近两万次。多次获得最佳论文奖、最高引用论文奖、Springer 十大畅销华人作者、Elsevier 最高引中国学者等。受邀担任了包括 SIGIR、WWW、KDD、ICML、NIPS、AAAI、ACL 在内的顶级国际会议的程序委员会主席或领域主席和多家国际学术期刊副主编。被聘为卡内基 - 梅隆大学 (CMU) 客座教授，诺丁汉大学荣誉教授，中国科技大学教授、博士生导师；被评为国际电子电气工程师学会 (IEEE) 会士，国际计算机学会 (ACM) 杰出会员。担任中国计算机学会青工委副主任，中文信息学会信息检索专委会副主任，中国云体系创新战略联盟常务理事。他的团队发布了 LightLDA、LightGBM、Multiverso 等知名的机器学习开源项目，并且为微软 CNTK 项目提供了分布式训练的解决方案，他的团队所参与的开源项目在 GitHub 上已累计获得数万颗星。



陈薇 微软亚洲研究院机器学习组主管研究员，研究机器学习各个分支的理论解释和算法改进，尤其关注深度学习、分布式机器学习、强化学习、博弈机器学习、排序学习等。2011 年于中国科学院数学与系统科学研究院获得博士学位，同年加入微软亚洲研究院，负责机器学习理论项目，先后在 NIPS、ICML、AAAI、IJCAI 等相关领域顶级国际会议和期刊上发表文章 30 余篇。



王太峰 蚂蚁金服人工智能部总监、资深算法专家。在蚂蚁金服负责 AI 算法组件建设，包括文本理解、图像理解、在线学习、强化学习等，算法工作服务于蚂蚁金服的支付、国际、保险等多条业务线。在加入蚂蚁之前在微软亚洲研究院工作 11 年，任主管研究员，他的研究方向包括大规模机器学习、数据挖掘、计算广告学等。在国际顶级的机器学习会议上发表近 20 篇论文，做了 4 次大规模机器学习专题讲座，并被多次邀请为各个会议程序委员。目前还是中国人工智能开源软件发展联盟的副秘书长，在大规模机器学习开源工具方面也做出了很多贡献，在微软期间主持开发过 DMTK 的开源项目，在 GitHub 上获得的点赞总数超过 8000 次，得到广泛好评。



高飞 微软亚洲研究院副研究员，主要从事分布式机器学习和深度学习的研究工作，并在国际会议上发表多篇论文。2014 年设计开发了当时规模最大的主题模型算法和系统 LightLDA。还开发了一系列分布式机器学习系统，并通过微软分布式机器学习工具包 (DMTK) 开源在 GitHub 上。



序言一	2.6.2 泛化误差的分解	/ 34
序言二	2.6.3 基于容量的估计误差的	
前 言	上界	/ 35
作者介绍	2.7 总结	/ 36
第1章 绪论	参考文献	/ 36
1.1 人工智能及其飞速发展		
1.2 大规模、分布式机器学习	第3章 分布式机器学习框架	/ 41
1.3 本书的安排	3.1 大数据与大模型的挑战	/ 42
参考文献	3.2 分布式机器学习的基本流程	/ 44
第2章 机器学习基础	3.3 数据与模型划分模块	/ 46
2.1 机器学习的基本概念	3.4 单机优化模块	/ 48
2.2 机器学习的基本流程	3.5 通信模块	/ 48
2.3 常用的损失函数	3.5.1 通信的内容	/ 48
2.3.1 Hinge 损失函数	3.5.2 通信的拓扑结构	/ 49
2.3.2 指数损失函数	3.5.3 通信的步调	/ 51
2.3.3 交叉熵损失函数	3.5.4 通信的频率	/ 52
2.4 常用的机器学习模型	3.6 数据与模型聚合模块	/ 53
2.4.1 线性模型	3.7 分布式机器学习理论	/ 54
2.4.2 核方法与支持向量机	3.8 分布式机器学习系统	/ 55
2.4.3 决策树与 Boosting	3.9 总结	/ 56
2.4.4 神经网络	参考文献	/ 57
2.5 常用的优化方法	第4章 单机优化之确定性算法	/ 61
2.6 机器学习理论	4.1 基本概述	/ 62
2.6.1 机器学习算法的泛化误差	4.1.1 机器学习的优化框架	/ 62

4.1.2	优化算法的分类和 发展历史	/ 65	5.3.4	等级优化算法	/ 107
4.2	一阶确定性算法	/ 67	5.4	总结	/ 109
4.2.1	梯度下降法	/ 67		参考文献	/ 109
4.2.2	投影次梯度下降法	/ 69	第6章 数据与模型并行		/ 113
4.2.3	近端梯度下降法	/ 70	6.1	基本概述	/ 114
4.2.4	Frank-Wolfe 算法	/ 71	6.2	计算并行模式	/ 117
4.2.5	Nesterov 加速法	/ 72	6.3	数据并行模式	/ 119
4.2.6	坐标下降法	/ 75	6.3.1	数据样本划分	/ 120
4.3	二阶确定性算法	/ 75	6.3.2	数据维度划分	/ 123
4.3.1	牛顿法	/ 76	6.4	模型并行模式	/ 123
4.3.2	拟牛顿法	/ 77	6.4.1	线性模型	/ 123
4.4	对偶方法	/ 78	6.4.2	神经网络	/ 127
4.5	总结	/ 81	6.5	总结	/ 133
	参考文献	/ 81		参考文献	/ 133
第5章 单机优化之随机算法		/ 85	第7章 通信机制		/ 135
5.1	基本随机优化算法	/ 86	7.1	基本概述	/ 136
5.1.1	随机梯度下降法	/ 86	7.2	通信的内容	/ 137
5.1.2	随机坐标下降法	/ 88	7.2.1	参数或参数的更新	/ 137
5.1.3	随机拟牛顿法	/ 91	7.2.2	计算的中间结果	/ 137
5.1.4	随机对偶坐标上升法	/ 93	7.2.3	讨论	/ 138
5.1.5	小结	/ 95	7.3	通信的拓扑结构	/ 139
5.2	随机优化算法的改进	/ 96	7.3.1	基于迭代式 MapReduce/ AllReduce 的通信拓扑	/ 140
5.2.1	方差缩减方法	/ 96	7.3.2	基于参数服务器的 通信拓扑	/ 142
5.2.2	算法组合方法	/ 100	7.3.3	基于数据流的通信拓扑	/ 143
5.3	非凸随机优化算法	/ 101	7.3.4	讨论	/ 145
5.3.1	Ada 系列算法	/ 102	7.4	通信的步调	/ 145
5.3.2	非凸理论分析	/ 104	7.4.1	同步通信	/ 146
5.3.3	逃离鞍点问题	/ 106			

7.4.2	异步通信	/ 147	9.3	异步算法	/ 187
7.4.3	同步和异步的平衡	/ 148	9.3.1	异步 SGD	/ 187
7.4.4	讨论	/ 150	9.3.2	Hogwild! 算法	/ 189
7.5	通信的频率	/ 150	9.3.3	Cyclades 算法	/ 190
7.5.1	时域滤波	/ 150	9.3.4	带延迟处理的异步算法	/ 192
7.5.2	空域滤波	/ 153	9.3.5	异步方法的进一步加速	/ 199
7.5.3	讨论	/ 155	9.3.6	讨论	/ 199
7.6	总结	/ 156	9.4	同步和异步的对比与融合	/ 199
参考文献	/ 156		9.4.1	同步和异步算法的 实验对比	/ 199
第 8 章 数据与模型聚合	/ 159		9.4.2	同步和异步的融合	/ 201
8.1	基本概念	/ 160	9.5	模型并行算法	/ 203
8.2	基于模型加和的聚合方法	/ 160	9.5.1	DistBelief	/ 203
8.2.1	基于全部模型加和的 聚合	/ 160	9.5.2	AlexNet	/ 204
8.2.2	基于部分模型加和的 聚合	/ 162	9.6	总结	/ 205
8.3	基于模型集成的聚合方法	/ 167	参考文献	/ 205	
8.3.1	基于输出加和的聚合	/ 168	第 10 章 分布式机器学习理论	/ 209	
8.3.2	基于投票的聚合	/ 171	10.1	基本概念	/ 210
8.4	总结	/ 174	10.2	收敛性分析	/ 210
参考文献	/ 174		10.2.1	优化目标和算法	/ 211
第 9 章 分布式机器学习算法	/ 177		10.2.2	数据和模型并行	/ 213
9.1	基本概念	/ 178	10.2.3	同步和异步	/ 215
9.2	同步算法	/ 179	10.3	加速比分析	/ 217
9.2.1	同步 SGD 方法	/ 179	10.3.1	从收敛速率到加速比	/ 218
9.2.2	模型平均方法及其改进	/ 182	10.3.2	通信量的下界	/ 219
9.2.3	ADMM 算法	/ 183	10.4	泛化分析	/ 221
9.2.4	弹性平均 SGD 算法	/ 185	10.4.1	优化的局限性	/ 222
9.2.5	讨论	/ 186	10.4.2	具有更好泛化能力的 非凸优化算法	/ 224

10.5 总结	/ 226	11.4 基于数据流的分布式机器学习系统	/ 241
参考文献	/ 226	11.4.1 数据流	/ 241
第 11 章 分布式机器学习系统	/ 229	11.4.2 TensorFlow 数据流系统	/ 243
11.1 基本概述	/ 230	11.5 实战比较	/ 248
11.2 基于 IMR 的分布式机器学习系统	/ 231	11.6 总结	/ 252
11.2.1 IMR 和 Spark	/ 231	参考文献	/ 252
11.2.2 Spark MLlib	/ 234	第 12 章 结语	/ 255
11.3 基于参数服务器的分布式机器学习系统	/ 236	12.1 全书总结	/ 256
11.3.1 参数服务器	/ 236	12.2 未来展望	/ 257
11.3.2 Multiverso 参数服务器	/ 237	索引	/ 260

CHAPTER

1

第1章

绪 论

- 1.1 人工智能及其飞速发展
- 1.2 大规模、分布式机器学习
- 1.3 本书的安排

DISTRIBUTED MACHINE LEARNING
Theories, Algorithms, and Systems

1.1 人工智能及其飞速发展

很早以前人类就有一个梦想：创建一种能像自己一样，具有独立思考和推理能力的机器。这个梦想驱动着人们不断进行科学探索，也孕育出很多引人入胜的科幻小说。

直到1956年一群怀揣梦想的青年科学家在美国的达特茅斯学院集会，正式提出了“人工智能”这一概念，从此开启了人工智能的历史篇章。在随后的60多年里，人工智能几起几落，技术也不断推陈出新，其波澜壮阔的景象如图1.1所示。每一次人工智能高潮（即所谓的“人工智能的春天”）的到来，都是因为某（几）项新技术的发明解决了之前困扰大家多年的难题，引燃了大众对于梦想的无限畅想和狂热追逐；而人工智能低谷（即所谓的“人工智能的冬天”）的出现，则往往是因为技术的发展速度跟不上狂热大众的期望，于是很多关于机器智能的预言破灭，政府和投资机构相继撤资，导致人工智能的研究得不到应有的充分支持。与其他的研究领域相比，正因为人工智能与人类本身更加密切相关，离人类的梦想更近，所以难免命运多舛，跌宕起伏。

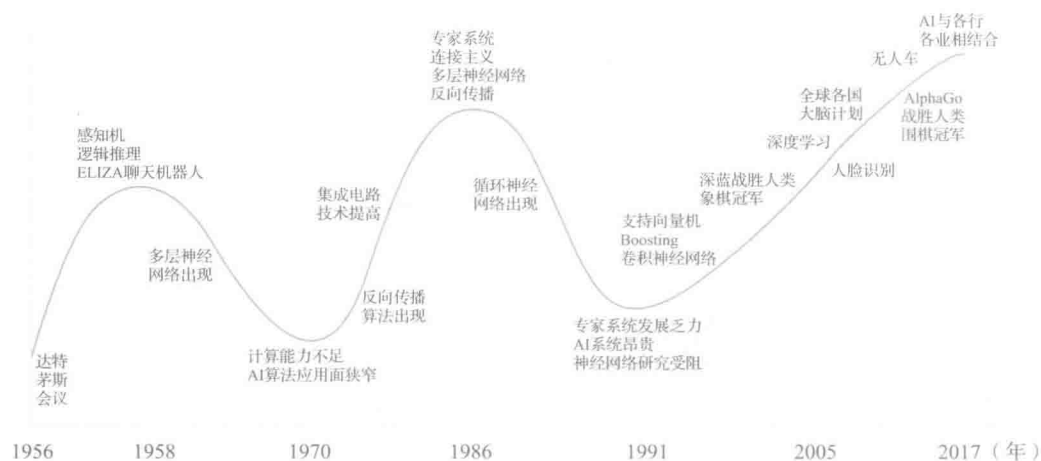


图 1.1 人工智能的发展历程

聚焦过去的十年间，人工智能技术取得了前所未有的高速发展，甚至用“春天”都不足以形容人工智能的热度。为了给大家一个具象化的感觉，我们列举了这几年人工智能在语音、图像、自然语言处理、人机对弈、自动驾驶、医疗健康等方面所取得的骄人战绩。

- 在语音处理方面，2016年年底，来自微软研究院的计算机科学家首次在普通对话数据上取得了可以和人类媲美的语音识别能力（词错误率低至5.9%）^[1]；2017年年初，IBM的科学家宣布通过集成多个语音识别模型，可以把识别的错误率进一步降低到5.1%^[2]；而2017年8月，微软研究院的科学家再接再厉，挑战技术极限，成功地训练出可以取得5.1%错误率的单个语音识别模型^[3]。语音识别精度的不断提升，为个人语音助手和智能音箱的出现提供了技术支持，催生出如Siri、Cortana、Google Now、Echo等家喻户晓的产品。
- 在图像处理方面，2012年AlexNet在ImageNet大规模视觉识别挑战（ILSVRC）中将图像分类的top-5错误率降低到15.3%^[4]；而后VGGNet将这个分类错误率进一步降低到7.3%^[5]；2015年，来自微软研究院和谷歌的科学家分别独立取得了错误率接近3.5%的骄人成绩^[6-7]，而这个识别精度已经远超人类的平均水平（5.1%）。图像识别能力的提升，促进了安防、智能金融等诸多领域的迅猛发展，造就了一批“刷脸公司”，不断刷新创业公司融资的纪录。
- 在自然语言处理方面，各大公司近年来都在构建自己的机器翻译系统，并且将其产品化、商业化。比如2016年年底微软公司发布的新版Microsoft Translator手机应用可以通过对100种语言的实时互译，实现跨国团队之间的无障碍实时交流。而谷歌、脸书、百度等公司也不断推出新型的机器翻译模型，在精度和速度上开展“军备竞赛”。2018年年初，微软公司宣布在中英新闻翻译领域达到了人类的水平，创立了人工智能的又一个里程碑^[8]。
- 在人机对弈方面，人工智能技术对人类选手构成了前所未有的威胁。2016年年初，来自DeepMind的AlphaGo以4:1的大比分战胜人类围棋世界冠军李世石^[9]。一年之后，新一代的AlphaGo化身Master，又在围棋快棋赛中横扫人类选手获得60连胜，随后又完胜当时世界排名第一的中国棋手柯洁^[10]。除了围棋，人工智能技术在德州扑克、桥牌、麻将等竞技领域也捷报频传^[11-13]。
- 在自动驾驶方面，从芯片制造商、互联网公司到传统汽车企业都开始引入人工智能技术，为未来布局。例如谷歌、百度、特斯拉、优步等公司先后宣布了自动驾驶的战略，而且在一定范围内实现了路测；而近期英特尔公司以153亿美元的天价收购了Mobileye公司，提升其在自动驾驶领域的核心竞争力^[14]。国内以自动驾驶为主题的创业公司更是如雨后春笋，关注自动驾驶的不同技术模块或不同场景，不断刺激着人们对自动驾驶走入寻常百姓家的美好憧憬。

- 在健康医疗方面，人工智能公司和传统医药企业密切合作，也取得了很多可喜的进展。例如，由谷歌的科学家训练出的人工智能模型在皮肤癌检测上达到了专业医师的水平^[15]；来自微软的科学家和医药公司及医院一起，建立医疗知识图谱和医师助手，并且利用深度学习技术把对糖尿病眼盲症的自动诊断引入临床^[16]。2017年年末，微软宣布了 Azure 云平台对基因数据分析的支持，包括大规模基因数据、基因分析 API 等，为进一步应用人工智能技术解决健康医疗的难题打下基础^[17-18]。

除了上述领域之外，人工智能在金融、物流、教育、制造等方面也都取得了长足的进步。这一次人工智能的热潮不仅仅是技术层面的，还涉及广泛的产业和资本运作。无论人们如何评价目前的产业状况，毋庸置疑的是，人工智能真的来了，而且即将对我们的生活产生巨大的影响。

1.2 大规模、分布式机器学习

人工智能真的来了。有人称 2016 年为人工智能元年、2017 年为人工智能的落地之年。大众对于人工智能的认知达到了前所未有的程度，传统产业对于智能转型的热情也空前高涨。那么，是什么原因导致人工智能的全面爆发呢？是在人工智能的算法和理论层面上出现了革命性的突破吗？反观当今主流的人工智能技术，我们会发现其实绝大部分机器学习算法（至少其原型）都是上世纪八九十年代，甚至更早就被提出来的，虽然近年来人们进行了很多技术改良，但是尚谈不上有革命性的技术突破。但是，在很多其他方面，却实实在在发生着改变，从最初的量变，逐步发展到质变，成为人工智能蓬勃发展的强力推手。用一个字来总结这种改变，应该就是“大”：在前所未有的大数据（尤其是有标签的训练数据）的支撑下，通过庞大的计算机集群（尤其以 GPU 集群为主），训练大规模的机器学习模型（尤其是深层神经网络）。而如此训练出来的机器学习模型因为足够复杂，可以有效地逼近很多困难问题的决策边界，因此可以最终秒杀传统的人工智能技术。

大数据、大模型为人工智能的飞速发展奠定了坚实的物质基础，也提出了新的技术挑战。近年来，越来越多的学者开始深入研究分布式机器学习，从而可以更高效地利用大数据训练更准确的大模型。分布式机器学习涉及如何分配训练任务，调配计算资源，协调各个功能模块，以达到训练速度与精度的平衡。一个分布式机器学习系统通常会包