

## MACHINE LEARNING ALGORITHM

深入分析机器学习中的常用算法，兼顾算法、理论与实践，帮助读者快速掌握算法精髓！

# Python 机器学习算法

赵志勇◎著

探索数据的内在价值，  
洞悉人工智能背后的技术！



中国工信出版集团



电子工业出版社  
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY  
<http://www.phei.com.cn>

# Python

## 机器学习算法



赵志勇◎著

电子工业出版社  
Publishing House of Electronics Industry  
北京•BEIJING

## 内 容 简 介

本书是一本机器学习入门读物，注重理论与实践的结合。全书主要包括 6 个部分，每个部分均以典型的机器学习算法为例，从算法原理出发，由浅入深，详细介绍算法的理论，并配合目前流行的 Python 语言，从零开始，实现每一个算法，以加强对机器学习算法理论的理解、增强实际的算法实践能力，最终达到熟练掌握每一个算法的目的。与其他机器学习类图书相比，本书同时包含算法理论的介绍和算法的实践，以理论支撑实践，同时，又将复杂、枯燥的理论用简单易懂的形式表达出来，促进对理论的理解。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

### 图书在版编目（CIP）数据

Python 机器学习算法 / 赵志勇著. —北京：电子工业出版社，2017.7

ISBN 978-7-121-31319-6

I .①P… II .①赵… III.①软件工具—程序设计 ②机器学习 IV.①TP311.561 ②TP181

中国版本图书馆 CIP 数据核字(2017)第 072742 号

策划编辑：符隆美

责任编辑：徐津平

印 刷：北京季蜂印刷有限公司

装 订：北京季蜂印刷有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本：720×1000 1/16 印张：22.75 字数：426 千字

版 次：2017 年 7 月第 1 版

印 次：2017 年 7 月第 1 次印刷

定 价：69.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，联系及邮购电话：(010) 88254888, 88258888。

质量投诉请发邮件至 [zlts@phei.com.cn](mailto:zlts@phei.com.cn)，盗版侵权举报请发邮件至 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

本书咨询联系方式：010-51260888-819, [faq@phei.com.cn](mailto:faq@phei.com.cn)。

# 推荐序

志勇是我在新浪微博的同事，刚来的时候坐我的旁边。记得当时志勇喜欢把看过的论文的重点部分剪下来粘到自己的笔记本上，并用五颜六色的笔标注，后来还知道志勇平时会写博客来记录自己在算法学习和实践中的心得。由此可见，志勇是一个非常认真且善于归纳总结的人。2016年志勇告诉我他在写这样一本书的时候，我深以为然，感觉正确的人做了一件正确的事。

志勇的书快完成的时候邀请我写序，这对我绝对是个挑战。幸运的是，书中的内容是我所熟悉的，仿佛是发生在自己身边的事。写作风格也和志勇平时交流时一致。所以，我可以从故事参与者的角度去介绍一下这本书。

说到机器学习算法，这两年可谓蓬勃发展。AlphaGo 战胜世界围棋冠军李世石已经成了大家茶余饭后的谈资，无人驾驶汽车是资本竞相追逐的万亿级市场，这些都源于数据收集能力、计算能力的提升，以及智能设备的普及。机器学习也已经在我们身边一些领域取得了成功，例如，现在已经获得上亿用户使用的今日头条。今日头条通过抓取众多媒体的资讯，利用机器学习算法推荐给用户，从而做到了资讯量大、更新快、更加个性化。

作为一个互联网从业者，更是感觉到机器学习算法已经融入到越来越多的产品和功能中。我曾经在一次交流中得知，某个应用为减少用户输入，用了一个团队的力量做输入选项的推荐。这在以前是不曾出现过的，以前的网站更多的是增加功能，让用户选择。现在更多的是推荐给用户，帮助用户选择。这里面既有移动应用屏幕小、操作复杂的原因，也有互联网公司越来越重视用户体验的原因。所以，在此要恭喜这本书的读者，你们选择了一个前途光明的行业。

机器学习算法比较典型的应用是推荐、广告和搜索。我们利用协同过滤技术来推荐商品、利用逻辑回归技术来做点击率的预测、利用分类技术来识别“垃圾”网页等。学好、用好每一种算法都很困难，需要掌握背后的理论基础，以及进行大量的实践，否则就会浮于表面，模仿他人，不能根据自己的业务做出合理的选择。

而市场上的书通常要么只是一些概要性质的介绍，要么是偏向实战，理论基础介绍得比较少。本书是少有的两者兼具的书，每一种算法都先介绍数学基础，再用 Python 代码做简单版本的实现，并且算法之间循序渐进，层层深入，读来如沐春风。

本书介绍了 LR、FM、SVM、协同过滤、矩阵分解等推荐和广告领域常用算法，有很强的实用性。深度学习更是近期主流互联网公司研究的热门领域。无论对机器学习的初学者还是已经具备一些项目经验的人来说，这都是很好的读本。希望本书对更多的人有益，也希望中国的“人工智能+”蓬勃发展。

新浪微博算法经理 陶辉

# 前言

## 起源

在读研究生期间，我就对机器学习算法萌生了很浓的兴趣，并对机器学习中的常用算法进行了学习，利用 MATLAB 对每一个算法进行了实践。在此过程中，每当遇到不懂的概念或者算法时，就会在网上查找相关的资料。也看到很多人在博客中分享算法的学习心得及算法的具体过程，其中有不少内容让我受益匪浅，但是有的内容仅仅是算法的描述，缺少实践的具体过程。

注意到这一点之后，我决定开始在博客中分享自己学习每一个机器学习算法的点点滴滴，为了让更多的初学者能够理解算法的具体过程并从中受益，我计划从三个方面出发，第一是算法过程的简单描述，第二是算法理论的详细推导，第三是算法的具体实践。2014 年 1 月 10 日，我在 CSDN 上写下了第一篇博客。当时涉及的方向主要是优化算法和简单易学的机器学习算法。

随着学习的深入，博客的内容越来越多，同时，在写作过程中，博客的质量也在慢慢提高，这期间也是机器学习快速发展的阶段，在行业内出现了很多优秀的算法库，如 Java 版本的 weka、Python 版本的 sklearn，以及其他的一些开源程序，通过对这些算法库的学习，我丰富了很多算法的知识，同时，我将学习到的心得记录在简单易学的机器学习算法中。工作之后，越发觉得这些基础知识对于算法的理解很有帮助，积累的这些算法学习材料成了我宝贵的财富。

2016 年，电子工业出版社博文视点的符隆美编辑联系到我，询问我是否有意向将这些博文汇总写一本书。能够写一本书是很多人的梦想，我也不例外。于是在 2016

年 9 月，我开始了对本书的构思，从选择算法开始，选择出使用较多的一些机器学习算法。在选择好算法后，从算法原理和算法实现两个方面对算法进行描述，希望本书能够在内容上既能照顾到初学者，又能使具有一定机器学习基础的读者从中受益。

在写作的过程中，我重新查阅了资料，力求保证知识的准确性，同时，在实践的环节中，我使用了目前比较流行的 Python 语言实现每一个算法，使得读者能够更容易理解算法的过程，在介绍深度学习的部分时，使用到了目前最热门的 TensorFlow 框架。为了帮助读者理解机器学习算法在实际工作中的具体应用，本书专门有一章介绍项目实践的部分，综合前面各种机器学习算法，介绍每一类算法在实际工作中的具体应用。

## 内容组织

本书开篇介绍机器学习的基本概念，包括监督学习、无监督学习和深度学习的基本概念。

第一部分介绍分类算法。分类算法是机器学习中最常用的算法。在分类算法中着重介绍 Logistic 回归、Softmax Regression、Factorization Machine、支持向量机、随机森林和 BP 神经网络等算法。

第二部分介绍回归算法。与分类算法不同的是，在回归算法中其目标值是连续的值，而在分类算法中，其目标值是离散的值。在回归算法中着重介绍线性回归、岭回归和 CART 树回归。

第三部分介绍聚类算法。聚类是将具有某种相同属性的数据聚成一个类别。在聚类算法中着重介绍 K-Means 算法、Mean Shift 算法、DBSCAN 算法和 Label Propagation 算法。

第四部分介绍推荐算法。推荐算法是一类基于具体应用场景的算法的总称。在推荐算法中着重介绍基于协同过滤的推荐、基于矩阵分解的推荐和基于图的推荐。

第五部分介绍深度学习。深度学习是近年来研究最为火热的方向。深度学习的网络模型和算法有很多种，在本书中，主要介绍最基本的两种算法：AutoEncoder 和卷积神经网络。

第六部分介绍以上这些算法在具体项目中的实践。通过具体的例子，可以清晰地看到每一类算法的应用场景。

附录介绍在实践中使用到的 Python 语言、numpy 库及 TensorFlow 框架的具体使用方法。

## 小结

本书试图从算法原理和实践两个方面来介绍机器学习中的常用算法，对每一类机器学习算法，精心挑选了具有代表性的算法，从理论出发，并配以详细的代码，本书的所有示例代码使用 Python 语言作为开发语言，读者可以从 <https://github.com/zhaozhiyong19890102/Python-Machine-Learning-Algorithm> 中下载本书的全部示例代码。

由于时间仓促，书中难免存在错误，欢迎广大读者和专家批评、指正，同时，欢迎大家提供意见和反馈。本书作者的电子邮箱：[zhaozhiyong1989@126.com](mailto:zhaozhiyong1989@126.com)。

## 致谢

首先，我要感谢陶辉和孙永生这两位良师益友，在本书的写作过程中，为我提供了很多意见和建议，包括全书的组织架构。感谢陶辉抽出宝贵的时间帮我写序，感谢孙永生帮我检查程序。

其次，我要感谢符隆美编辑和董雪编辑在写作和审稿的过程中对我的鼓励和悉心指导。

再次，我要感谢姜贵彬、易慧民、潘文彬，感谢他们能够抽出宝贵的时间帮本书写推荐语，感谢他们在读完本书后给出的宝贵意见和建议。

然后，我要感谢 July 在本书的写作过程中对本书提出的宝贵意见，感谢张俊林、王斌在读完本书初稿后对本书的指点。

最后，感谢我的亲人和朋友，是你们的鼓励才使得本书能够顺利完成。

赵志勇

2017 年 6 月 6 日于北京

## 读者服务

轻松注册成为博文视点社区用户 ([www.broadview.com.cn](http://www.broadview.com.cn))，扫码直达本书页面。

- **提交勘误：**您对书中内容的修改意见可在 提交勘误 处提交，若被采纳，将获赠博文视点社区积分（在您购买电子书时，积分可用来抵扣相应金额）。
- **交流互动：**在页面下方 读者评论 处留下您的疑问或观点，与我们和其他读者一同学习交流。

页面入口：<http://www.broadview.com.cn/31319>



三步加入“人工智能交流群”，实时获取资源共享，并有机会与大咖实时交流。

- 扫码添加小编为微信好友。
- 申请验证时输入“AI”。
- 小编带你加入“人工智能交流群”。



# 目录

0 绪论 .....	1
0.1 机器学习基础 .....	1
0.1.1 机器学习的概念 .....	1
0.1.2 机器学习算法的分类 .....	2
0.2 监督学习 .....	3
0.2.1 监督学习 .....	3
0.2.2 监督学习的流程 .....	3
0.2.3 监督学习算法 .....	4
0.3 无监督学习 .....	4
0.3.1 无监督学习 .....	4
0.3.2 无监督学习的流程 .....	4
0.3.3 无监督学习算法 .....	5
0.4 推荐系统和深度学习 .....	6
0.4.1 推荐系统 .....	6
0.4.2 深度学习 .....	6
0.5 Python 和机器学习算法实践 .....	6
参考文献 .....	7

## 第一部分 分类算法

1 Logistic Regression .....	10
1.1 Logistic Regression 模型 .....	10

1.1.1 线性可分 VS 线性不可分 .....	10
1.1.2 Logistic Regression 模型 .....	11
1.1.3 损失函数 .....	13
1.2 梯度下降法 .....	14
1.2.1 梯度下降法的流程 .....	14
1.2.2 凸优化与非凸优化 .....	15
1.2.3 利用梯度下降法训练 Logistic Regression 模型 .....	17
1.3 梯度下降法的若干问题 .....	18
1.3.1 选择下降的方向 .....	18
1.3.2 步长的选择 .....	19
1.4 Logistic Regression 算法实践 .....	20
1.4.1 利用训练样本训练 Logistic Regression 模型 .....	20
1.4.2 最终的训练效果 .....	22
1.4.3 对新数据进行预测 .....	23
参考文献 .....	26
<b>2 Softmax Regression .....</b>	<b>27</b>
2.1 多分类问题 .....	27
2.2 Softmax Regression 算法模型 .....	28
2.2.1 Softmax Regression 模型 .....	28
2.2.2 Softmax Regression 算法的代价函数 .....	28
2.3 Softmax Regression 算法的求解 .....	29
2.4 Softmax Regression 与 Logistic Regression 的关系 .....	31
2.4.1 Softmax Regression 中的参数特点 .....	31
2.4.2 由 Softmax Regression 到 Logistic Regression .....	31
2.5 Softmax Regression 算法实践 .....	32
2.5.1 对 Softmax Regression 算法的模型进行训练 .....	33
2.5.2 最终的模型 .....	34
2.5.3 对新的数据的预测 .....	35
参考文献 .....	39
<b>3 Factorization Machine .....</b>	<b>40</b>
3.1 Logistic Regression 算法的不足 .....	40
3.2 因子分解机 FM 的模型 .....	42

3.2.1 因子分解机 FM 模型 .....	42
3.2.2 因子分解机 FM 可以处理的问题 .....	43
3.2.3 二分类因子分解机 FM 算法的损失函数 .....	43
3.3 FM 算法中交叉项的处理 .....	43
3.3.1 交叉项系数 .....	43
3.3.2 模型的求解 .....	44
3.4 FM 算法的求解 .....	45
3.4.1 随机梯度下降 (Stochastic Gradient Descent) .....	45
3.4.2 基于随机梯度的方式求解 .....	45
3.4.3 FM 算法流程 .....	46
3.5 因子分解机 FM 算法实践 .....	49
3.5.1 训练 FM 模型 .....	50
3.5.2 最终的训练效果 .....	53
3.5.3 对新的数据进行预测 .....	55
参考文献 .....	57
<b>4 支持向量机 .....</b>	<b>58</b>
4.1 二分类问题 .....	58
4.1.1 二分类的分隔超平面 .....	58
4.1.2 感知机算法 .....	59
4.1.3 感知机算法存在的问题 .....	61
4.2 函数间隔和几何间隔 .....	61
4.2.1 函数间隔 .....	62
4.2.2 几何间隔 .....	62
4.3 支持向量机 .....	63
4.3.1 间隔最大化 .....	63
4.3.2 支持向量和间隔边界 .....	64
4.3.3 线性支持向量机 .....	65
4.4 支持向量机的训练 .....	66
4.4.1 学习的对偶算法 .....	66
4.4.2 由线性支持向量机到非线性支持向量机 .....	68
4.4.3 序列最小最优化算法 SMO .....	69
4.5 支持向量机 SVM 算法实践 .....	74
4.5.1 训练 SVM 模型 .....	74

4.5.2 利用训练样本训练 SVM 模型.....	81
4.5.3 利用训练好的 SVM 模型对新数据进行预测.....	85
参考文献.....	88
<b>5 随机森林 .....</b>	<b>89</b>
5.1 决策树分类器 .....	89
5.1.1 决策树的基本概念 .....	89
5.1.2 选择最佳划分的标准 .....	91
5.1.3 停止划分的标准 .....	94
5.2 CART 分类树算法.....	95
5.2.1 CART 分类树算法的基本原理.....	95
5.2.2 CART 分类树的构建.....	95
5.2.3 利用构建好的分类树进行预测 .....	98
5.3 集成学习（Ensemble Learning） .....	99
5.3.1 集成学习的思想 .....	99
5.3.2 集成学习中的典型方法 .....	99
5.4 随机森林（Random Forests） .....	101
5.4.1 随机森林算法模型 .....	101
5.4.2 随机森林算法流程 .....	102
5.5 随机森林 RF 算法实践 .....	104
5.5.1 训练随机森林模型 .....	105
5.5.2 最终的训练结果 .....	109
5.5.3 对新数据的预测 .....	110
参考文献 .....	113
<b>6 BP 神经网络 .....</b>	<b>114</b>
6.1 神经元概述 .....	114
6.1.1 神经元的基本结构 .....	114
6.1.2 激活函数 .....	115
6.2 神经网络模型 .....	116
6.2.1 神经网络的结构 .....	116
6.2.2 神经网络中的参数说明 .....	117
6.2.3 神经网络的计算 .....	117
6.3 神经网络中参数的求解 .....	118

6.3.1 神经网络损失函数 .....	118
6.3.2 损失函数的求解 .....	119
6.3.3 BP 神经网络的学习过程 .....	120
6.4 BP 神经网络中参数的设置 .....	126
6.4.1 非线性变换 .....	126
6.4.2 权重向量的初始化 .....	126
6.4.3 学习率 .....	127
6.4.4 隐含层节点的个数 .....	127
6.5 BP 神经网络算法实践 .....	127
6.5.1 训练 BP 神经网络模型 .....	128
6.5.2 最终的训练效果 .....	132
6.5.3 对新数据的预测 .....	133
参考文献 .....	136

## 第二部分 回归算法

7 线性回归 .....	138
7.1 基本线性回归 .....	138
7.1.1 线性回归的模型 .....	138
7.1.2 线性回归模型的损失函数 .....	139
7.2 线性回归的最小二乘解法 .....	140
7.2.1 线性回归的最小二乘解法 .....	140
7.2.2 广义逆的概念 .....	141
7.3 牛顿法 .....	141
7.3.1 基本牛顿法的原理 .....	141
7.3.2 基本牛顿法的流程 .....	142
7.3.3 全局牛顿法 .....	142
7.3.4 Armijo 搜索 .....	144
7.3.5 利用全局牛顿法求解线性回归模型 .....	145
7.4 利用线性回归进行预测 .....	146
7.4.1 训练线性回归模型 .....	147
7.4.2 最终的训练结果 .....	149
7.4.3 对新数据的预测 .....	150
7.5 局部加权线性回归 .....	152

7.5.1 局部加权线性回归模型 .....	152
7.5.2 局部加权线性回归的最终结果 .....	153
参考文献 .....	154
<b>8 岭回归和 Lasso 回归 .....</b>	<b>155</b>
8.1 线性回归存在的问题 .....	155
8.2 岭回归模型 .....	156
8.2.1 岭回归模型 .....	156
8.2.2 岭回归模型的求解 .....	156
8.3 Lasso 回归模型 .....	157
8.4 拟牛顿法 .....	158
8.4.1 拟牛顿法 .....	158
8.4.2 BFGS 校正公式的推导 .....	158
8.4.3 BFGS 校正的算法流程 .....	159
8.5 L-BFGS 求解岭回归模型 .....	162
8.5.1 BFGS 算法存在的问题 .....	162
8.5.2 L-BFGS 算法思路 .....	162
8.6 岭回归对数据的预测 .....	165
8.6.1 训练岭回归模型 .....	166
8.6.2 最终的训练结果 .....	168
8.6.3 利用岭回归模型预测新的数据 .....	168
参考文献 .....	171
<b>9 CART 树回归 .....</b>	<b>172</b>
9.1 复杂的回归问题 .....	172
9.1.1 线性回归模型 .....	172
9.1.2 局部加权线性回归 .....	173
9.1.3 CART 算法 .....	174
9.2 CART 回归树生成 .....	175
9.2.1 CART 回归树的划分 .....	175
9.2.2 CART 回归树的构建 .....	177
9.3 CART 回归树剪枝 .....	179
9.3.1 前剪枝 .....	179
9.3.2 后剪枝 .....	180

9.4	CART 回归树对数据预测 .....	180
9.4.1	利用训练数据训练 CART 回归树模型 .....	180
9.4.2	最终的训练结果 .....	182
9.4.3	利用训练好的 CART 回归树模型对新的数据预测 .....	185
	参考文献 .....	187
 第三部分 聚类算法		
<b>10</b>	<b>K-Means .....</b>	<b>190</b>
10.1	相似性的度量 .....	190
10.1.1	闵可夫斯基距离 .....	191
10.1.2	曼哈顿距离 .....	191
10.1.3	欧氏距离 .....	191
10.2	K-Means 算法原理 .....	192
10.2.1	K-Means 算法的基本原理 .....	192
10.2.2	K-Means 算法步骤 .....	193
10.2.3	K-Means 算法与矩阵分解 .....	193
10.3	K-Means 算法实践 .....	195
10.3.1	导入数据 .....	196
10.3.2	初始化聚类中心 .....	197
10.3.3	聚类过程 .....	198
10.3.4	最终的聚类结果 .....	199
10.4	K-Means++算法 .....	200
10.4.1	K-Means 算法存在的问题 .....	200
10.4.2	K-Means++算法的基本思路 .....	202
10.4.3	K-Means++算法的过程和最终效果 .....	204
	参考文献 .....	205
<b>11</b>	<b>Mean Shift .....</b>	<b>206</b>
11.1	Mean Shift 向量 .....	206
11.2	核函数 .....	207
11.3	Mean Shift 算法原理 .....	209
11.3.1	引入核函数的 Mean Shift 向量 .....	209
11.3.2	Mean Shift 算法的基本原理 .....	210

11.4 Mean Shift 算法的解释 .....	212
11.4.1 概率密度梯度 .....	212
11.4.2 Mean Shift 向量的修正 .....	213
11.4.3 Mean Shift 算法流程 .....	213
11.5 Mean Shift 算法实践 .....	217
11.5.1 Mean Shift 的主过程 .....	218
11.5.2 Mean Shift 的最终聚类结果 .....	219
参考文献 .....	221
<b>12 DBSCAN .....</b>	<b>222</b>
12.1 基于密度的聚类 .....	222
12.1.1 基于距离的聚类算法存在的问题 .....	222
12.1.2 基于密度的聚类算法 .....	225
12.2 DBSCAN 算法原理 .....	225
12.2.1 DBSCAN 算法的基本概念 .....	225
12.2.2 DBSCAN 算法原理 .....	227
12.2.3 DBSCAN 算法流程 .....	228
12.3 DBSCAN 算法实践 .....	231
12.3.1 DBSCAN 算法的主要过程 .....	232
12.3.2 Mean Shift 的最终聚类结果 .....	234
参考文献 .....	236
<b>13 Label Propagation .....</b>	<b>237</b>
13.1 社区划分 .....	237
13.1.1 社区以及社区划分 .....	237
13.1.2 社区划分的算法 .....	238
13.1.3 社区划分的评价标准 .....	239
13.2 Label Propagation 算法原理 .....	239
13.2.1 Label Propagation 算法的基本原理 .....	239
13.2.2 标签传播 .....	240
13.2.3 迭代的终止条件 .....	242
13.3 Label Propagation 算法过程 .....	244
13.4 Label Propagation 算法实践 .....	244
13.4.1 导入数据 .....	245