

人气视频课程讲师

累计超过 30 万学员的选择

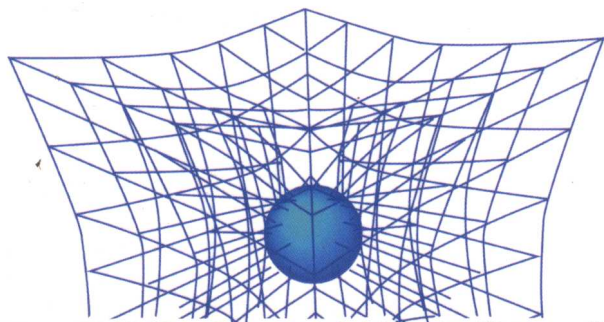


跟着迪哥学

Python

数据分析与机器学习实战

唐宇迪 / 著



算法原理 + 数学推导 + 项目实战 + 源代码 + 在线服务
机器学习实战课程套餐

面向零基础，通俗易懂，带给你沉浸式学习体验。

打通从算法原理、数学推导到实例操作的疑难点，快速入门人工智能领域。



中国工信出版集团



人民邮电出版社
POSTS & TELECOM PRESS

人气视频课程讲师



累计超过 30 万学员的选择

跟着迪哥学

Python

数据分析与机器学习实战

唐宇迪 / 著

人民邮电出版社

北京

图书在版编目(CIP)数据

跟着迪哥学Python数据分析与机器学习实战 / 唐宇迪著. — 北京: 人民邮电出版社, 2019.9
ISBN 978-7-115-51244-4

I. ①跟… II. ①唐… III. ①软件工具—程序设计
IV. ①TP311.561

中国版本图书馆CIP数据核字(2019)第087638号

内 容 提 要

本书结合了机器学习、数据分析和 Python 语言, 通过案例以通俗易懂的方式讲解了如何将算法应用到实际任务。

全书共 20 章, 大致分为 4 个部分。第 1 部分介绍了 Python 必备的工具包, 包括科学计算库 Numpy、数据分析库 Pandas、可视化库 Matplotlib; 第 2 部分讲解了机器学习中的经典算法, 例如回归算法、决策树、集成算法、支持向量机、聚类算法等; 第 3 部分介绍了深度学习中的常用算法, 包括神经网络、卷积神经网络、递归神经网络; 第 4 部分是项目实战, 基于真实数据集, 将算法模型应用到实际业务中。

本书适合对人工智能、机器学习、数据分析等方向感兴趣的初学者和爱好者。

-
- ◆ 著 唐宇迪
责任编辑 俞彬
责任印制 马振武
 - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路 11 号
邮编 100164 电子邮件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
北京市艺辉印刷有限公司印刷
 - ◆ 开本: 800×1000 1/16
印张: 28.75
字数: 765 千字 2019 年 9 月第 1 版
印数: 1—3 000 册 2019 年 9 月北京第 1 次印刷
-

定价: 89.00 元

读者服务热线: (010) 81055256 印装质量热线: (010) 81055316

反盗版热线: (010) 81055315

广告经营许可证: 京东工商广登字 20170147 号

前言

PREFACE

人工智能的飞速发展，带来了丰富的机遇与挑战。机器学习算法工程师、数据挖掘工程师、大数据工程师等岗位的薪资在 IT 行业也颇丰。面对高薪与前沿技术的诱惑，越来越多的大学毕业生准备投身其中，但苦于缺乏指导性教材进行系统学习，非科班出身的大学毕业生更是缺乏相关数学基础。

很多同学认识我是通过在线课程或线下培训，机器学习培训工作已经伴我走过了近 4 个年头。在这期间，开发的线上就业课程 40 余门，参与的学员累计超过 30 万人，顺利完成企业与高校讲师培训 30 余场，直播课程百余场。忙碌之余，最大的收获就是收到同学们晒出的各大企业的 offer 与认可。

在培训工作中，同学们给我最多的反馈就是虽然能参考的资料有很多，但是都很难理解，尤其对于初学者而言，看各种公式就要晕掉了。这几年我也一直在思考如何讲解才能让大家更深刻、更轻松地了解机器学习中的每一个算法。

本书是我多年培训教学和学习心得的总结，最大的特色就是以接地气的方式向大家通俗地讲解算法原理与应用方法，让读者能够更轻松地去理解其中每一个复杂的算法。学习的目的肯定要在实际任务中发挥作用，我写作的初衷也是希望更多读者能将理论与实战方法应用到自己的业务中，所以本书整体风格是以实战为主，通过案例来解读如何将机器学习应用在实际的数据挖掘任务中。

本书面向的读者

本书主要面向对人工智能、机器学习、数据分析等方面有强烈兴趣的初学者和爱好者，通过本书的学习，读者能够掌握机器学习中经典算法原理推导、整体流程以及其中数学公式与各种参数的作用。案例全部采用当下流行的 Python 语言，从最基础的工具包开始讲起，让大家熟练使用 Python 及其数据科学工具包进行机器学习和数据挖掘领域的项目实战任务，并处理其中遇到的种种问题。

路线图

本书内容大体可以分为以下 4 个部分。

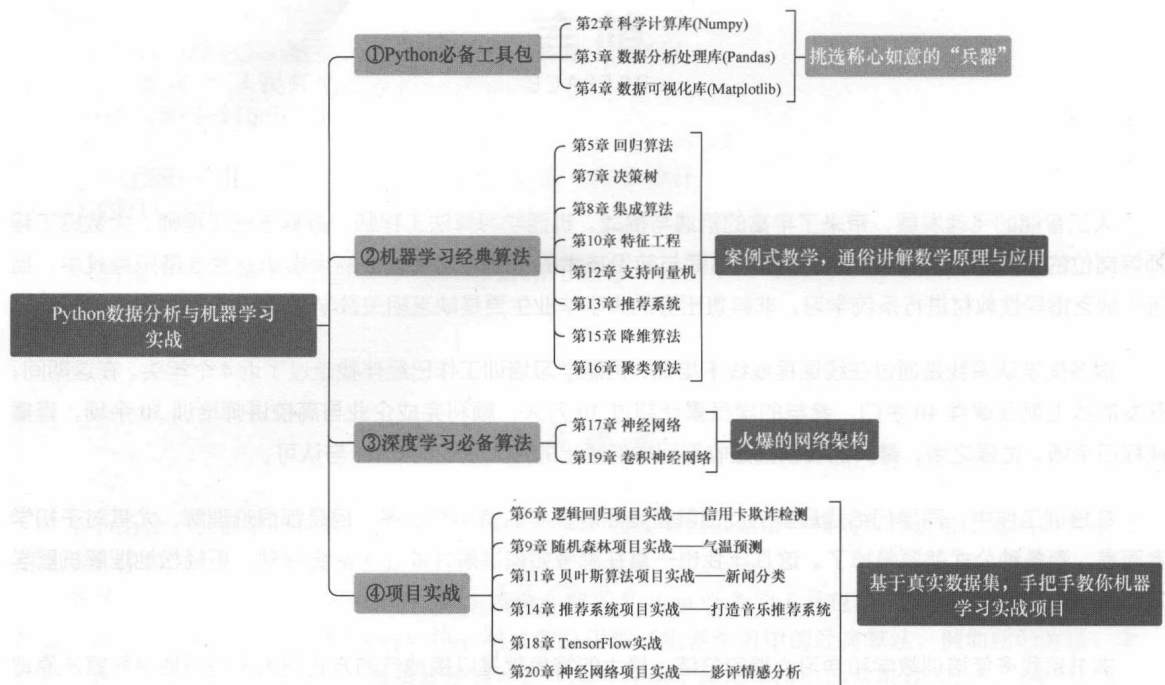


图 0-1 本书学习路线图

总结起来比较合适的学习路线如下。

第①步：Python 工具包的使用，先把称心如意的“兵器”准备好，它们是实战中的好帮手。

第②步：理解机器学习算法，建模分析的核心就是其中的算法了，打牢基础才能走得更远。

第③步：项目实战应用，将算法模型应用到实际业务中，通过实际任务来进行提升。

可能很多读者都觉得应当先把 Python 的基础打牢固再进行后续的学习，我觉得这样可能会花费较多时间，从而耽搁后续重点内容学习，建议读者对于编程语言通过实际案例边练边学，把重点放在机器学习原理与应用中。

阅读本书需要准备什么 / 如何使用本书

对于初学者来说，可能在学习路线以及职业规划上有些迷茫，这里结合我对机器学习与数据科学领域的理解来进行阐述分析。首先无论从事人工智能中哪个方向，肯定要从工程师做起，那手里一定得有一个称心如意的“兵器”，本书选择的是 Python 语言，基于 3.x 版本进行实战演示。读者如果具备大学数学基础，学习起来会相对更容易一些，在学习过程中，难免遇到各种难以理解的算法问题，建议大家先对其整体流程进

行通俗理解，再结合实际案例进行思考，很多时候数学上的描述十分复杂，而代码中的解释却浅显易懂。项目实战的目的的一方面是从应用的角度阐述如何进行实际任务建模与分析，另一方面也是一个积累的过程。人工智能行业发展迅速，不要停下学习的脚步，每天都要学习新的知识来充实自己。

配套资源

本书由异步社区 (<https://www.epubit.com/>) 为您提供相关资源。

本书提供配套的源代码和数据源文件。要想获得配套资源，请登陆异步社区，按书名搜索，进入本书页面，点击配套资源，跳转到下载页面，按提示进行操作即可。注意：为保证购书读者的权益，该操作会给出相关提示，要求输入提取码进行验证。

建议与反馈

由于作者水平有限，书中难免有错误和不当之处，欢迎读者指正。如果读者遇到问题需要帮助，也欢迎交流（微信号：digexiaozhushou），我期望与你共同成长。

目录

CONTENTS

第1章 人工智能入门指南	001	2.4.5 随机模块	031
1.1 AI时代首选 Python	002	2.4.6 文件读写	033
1.1.1 Python 的特点	002	本章总结	037
1.1.2 Python 该怎么学	002	第3章 数据分析处理库	038
1.2 人工智能的核心——机器学习	003	(Pandas)	038
1.2.1 什么是机器学习	003	3.1 数据预处理	039
1.2.2 机器学习的流程	004	3.1.1 数据读取	039
1.2.3 机器学习该怎么学	005	3.1.2 DataFrame 结构	040
1.3 环境配置	005	3.1.3 数据索引	042
1.3.1 Anaconda 大礼包	006	3.1.4 创建 DataFrame	046
1.3.2 Jupyter Notebook	009	3.1.5 Series 操作	048
1.3.3 上哪儿找资源	011	3.2 数据分析	051
本章总结	012	3.2.1 统计分析	051
第2章 科学计算库	013	3.2.2 pivot 数据透视表	055
(Numpy)	013	3.2.3 groupby 操作	058
2.1 Numpy 的基本操作	014	3.3 常用函数操作	063
2.1.1 array 数组	014	3.3.1 Merge 操作	063
2.1.2 数组特性	015	3.3.2 排序操作	066
2.1.3 数组属性操作	016	3.3.3 缺失值处理	067
2.2 索引与切片	017	3.3.4 apply 自定义函数	070
2.2.1 数值索引	017	3.3.5 时间操作	073
2.2.2 bool 索引	018	3.3.6 绘图操作	076
2.3 数据类型与数值计算	020	3.4 大数据处理技巧	079
2.3.1 数据类型	020	3.4.1 数值类型转换	079
2.3.2 复制与赋值	020	3.4.2 属性类型转换	082
2.3.3 数值运算	021	本章总结	084
2.3.4 矩阵乘法	024	第4章 数据可视化库	085
2.4 常用功能模块	025	(Matplotlib)	085
2.4.1 排序操作	025	4.1 常规绘图方法	086
2.4.2 数组形状操作	026	4.1.1 细节设置	086
2.4.3 数组的拼接	028	4.1.2 子图与标注	090
2.4.4 创建数组函数	029	4.1.3 风格设置	097

4.2 常用图表绘制	099	6.3.3 分类阈值对结果的影响	147
4.2.1 条形图	099	6.4 过采样方案	149
4.2.2 盒图	102	6.4.1 SMOTE 数据生成策略	150
4.2.3 直方图与散点图	105	6.4.2 过采样应用效果	151
4.2.4 3D 图	107	项目总结	152
4.2.5 布局设置	110	第 7 章 决策树	154
本章总结	111	7.1 决策树原理	155
第 5 章 回归算法	112	7.1.1 决策树的基本概念	155
5.1 线性回归算法	113	7.1.2 衡量标准	156
5.1.1 线性回归方程	113	7.1.3 信息增益	158
5.1.2 误差项分析	114	7.1.4 决策树构造实例	159
5.1.3 似然函数求解	115	7.1.5 连续值问题	161
5.1.4 线性回归求解	117	7.1.6 信息增益率	161
5.2 梯度下降算法	117	7.1.7 回归问题求解	162
5.2.1 下山方向选择	118	7.2 决策树剪枝策略	162
5.2.2 梯度下降优化	119	7.2.1 剪枝策略	162
5.2.3 梯度下降策略对比	120	7.2.2 决策树算法涉及参数	163
5.2.4 学习率对结果的影响	121	本章总结	164
5.3 逻辑回归算法	122	第 8 章 集成算法	165
5.3.1 原理推导	122	8.1 bagging 算法	166
5.3.2 逻辑回归求解	124	8.1.1 并行的集成	166
本章总结	125	8.1.2 随机森林	166
第 6 章 逻辑回归项目实战—— 信用卡欺诈检测	126	8.2 boosting 算法	170
6.1 数据分析与预处理	127	8.2.1 串行的集成	170
6.1.1 数据读取与分析	127	8.2.2 Adaboost 算法	171
6.1.2 样本不均衡解决方案	129	8.3 stacking 模型	173
6.1.3 特征标准化	129	本章总结	174
6.2 下采样方案	133	第 9 章 随机森林项目实战—— 气温预测	175
6.2.1 交叉验证	134	9.1 随机森林建模	176
6.2.2 模型评估方法	137	9.1.1 特征可视化与预处理	177
6.2.3 正则化惩罚	139	9.1.2 随机森林回归模型	183
6.3 逻辑回归模型	141	9.1.3 树模型可视化方法	184
6.3.1 参数对结果的影响	141		
6.3.2 混淆矩阵	144		

9.1.4 特征重要性	189	12.1.2 距离与标签定义	262
9.2 数据与特征对结果影响分析	192	12.1.3 目标函数	263
9.2.1 特征工程	194	12.1.4 拉格朗日乘法	264
9.2.2 数据量对结果影响分析	196	12.2 支持向量的作用	266
9.2.3 特征数量对结果影响分析	199	12.2.1 支持向量机求解	266
9.3 模型调参	206	12.2.2 支持向量的作用	267
9.3.1 随机参数选择	208	12.3 支持向量机涉及参数	268
9.3.2 网络参数搜索	212	12.3.1 软间隔参数的选择	268
项目总结	216	12.3.2 核函数的作用	270
第 10 章 特征工程	217	12.4 案例：参数对结果的影响	272
10.1 数值特征	218	12.4.1 SVM 基本模型	272
10.1.1 字符串编码	218	12.4.2 核函数变换	277
10.1.2 二值与多项式特征	222	12.4.3 SVM 参数选择	279
10.1.3 连续值离散化	225	12.4.4 SVM 人脸识别实例	281
10.1.4 对数与时间变换	228	本章总结	284
10.2 文本特征	230	第 13 章 推荐系统	285
10.2.1 词袋模型	230	13.1 推荐系统的应用	286
10.2.2 常用文本特征构造方法	234	13.2 协同过滤算法	288
10.3 论文与 benchmark	237	13.2.1 基于用户的协同过滤	288
本章总结	240	13.2.2 基于商品的协同过滤	291
第 11 章 贝叶斯算法项目实战—— 新闻分类	241	13.3 隐语义模型	292
11.1 贝叶斯算法	242	13.3.1 矩阵分解思想	292
11.1.1 贝叶斯公式	242	13.3.2 隐语义模型求解	294
11.1.2 拼写纠错实例	244	13.3.3 评估方法	296
11.1.3 垃圾邮件分类	246	本章总结	296
11.2 新闻分类任务	248	第 14 章 推荐系统项目实战—— 打造音乐推荐系统	297
11.2.1 数据清洗	249	14.1 数据集清洗	298
11.2.2 TF-IDF 关键词提取	253	14.1.1 统计分析	299
项目总结	259	14.1.2 数据集整合	303
第 12 章 支持向量机	260	14.2 基于相似度的推荐	308
12.1 支持向量机工作原理	261	14.2.1 排行榜推荐	309
12.1.1 支持向量机要解决的问题	261	14.2.2 基于歌曲相似度的推荐	310
		14.3 基于矩阵分解的推荐	313
		14.3.1 奇异值分解	313

14.3.2 使用 SVD 算法进行音乐推荐	317	17.2.1 整体框架	374
项目总结	322	17.2.2 神经元的作用	376
第 15 章 降维算法	323	17.2.3 正则化	378
15.1 线性判别分析	324	17.2.4 激活函数	379
15.1.1 降维原理概述	324	17.3 网络调优细节	381
15.1.2 优化的目标	325	17.3.1 数据预处理	381
15.1.3 线性判别分析求解	326	17.3.2 Drop-Out	382
15.1.4 Python 实现线性判别分析降维	328	17.3.3 数据增强	383
15.2 主成分分析	335	17.3.4 网络结构设计	384
15.2.1 PCA 降维基本知识	335	本章总结	384
15.2.2 PCA 优化目标求解	336	第 18 章 TensorFlow 实战	386
15.2.3 Python 实现 PCA 降维	338	18.1 TensorFlow 基本操作	387
本章总结	345	18.1.1 TensorFlow 特性	387
第 16 章 聚类算法	346	18.1.2 TensorFlow 基本操作	389
16.1 K-means 算法	347	18.1.3 TensorFlow 实现回归任务	392
16.1.1 聚类的基本特性	347	18.2 搭建神经网络进行手写字体识别	395
16.1.2 K-means 算法原理	348	本章总结	402
16.1.3 K-means 涉及参数	350	第 19 章 卷积神经网络	403
16.1.4 K-means 聚类效果与优缺点	352	19.1 卷积操作原理	404
16.2 DBSCAN 聚类算法	353	19.1.1 卷积神经网络应用	404
16.2.1 DBSCAN 算法概述	353	19.1.2 卷积操作流程	406
16.2.2 DBSCAN 工作流程	354	19.1.3 卷积计算方法	408
16.2.3 半径对结果的影响	357	19.1.4 卷积涉及参数	411
16.3 聚类实例	358	19.1.5 池化层	415
本章总结	363	19.2 经典网络架构	416
第 17 章 神经网络	364	19.2.1 卷积神经网络整体架构	416
17.1 神经网络必备基础	365	19.2.2 AlexNet 网络	417
17.1.1 神经网络概述	365	19.2.3 VGG 网络	418
17.1.2 计算机眼中的图像	367	19.2.4 ResNet 网络	421
17.1.3 得分函数	368	19.3 TensorFlow 实战卷积神经网络	424
17.1.4 损失函数	370	本章总结	427
17.1.5 反向传播	372		
17.2 神经网络整体架构	374		

第 20 章 神经网络项目实战—— 影评情感分析 428

20.1 递归神经网络 429

20.1.1 RNN 网络架构 429

20.1.2 LSTM 网络 430

20.2 影评数据特征工程 431

20.2.1 词向量 432

20.2.2 数据特征制作 436

20.3 构建 RNN 模型 444

项目总结 449

第 1 章

人工智能入门指南

当今时代，人工智能迅速发展，高薪的诱惑、前沿的技术挑战使得越来越多的小伙伴想要学习人工智能，那么更大的问题也就随之产生了——如何学习人工智能呢？正所谓“万事开头难”，如何走好第一步十分关键。学习人工智能的成本还是蛮高的，一般来说，付出了大量的时间和精力，一定要有满意的收获才可以。作为 Python 开篇之讲，本章首先介绍机器学习处理问题的方法与流程，以及实战必备武器——Python 基础教程及其环境配置。

1.1 AI 时代首选 Python

人工智能就是用编程实现各种算法和数据建模。提起编程，以前大家可能更注重 C 语言和 Java 语言，但是现在，Python 在数据科学领域运用广泛，相信大家早已在各大媒体和圈子中看到 Python 与日俱增的发展前景，可以说，Python 已经成为当下最火的编程语言之一了（见图 1-1）。



图 1-1 AI 时代首选 Python

1.1.1 Python 的特点

Python 被当作“核心武器”肯定是有原因的，进入 AI 行业，大家最初给自己的定位基本都是工程师，办事效率肯定是越高越好，这跟 Python 的出发点也是一致的，试问：能用 1 行代码解决的问题，何必用 10 行呢？

如果大家学过 C 语言，肯定会觉得它用起来还是比较麻烦的，限制非常多。但是用 Python 写起程序来可以更随性一些，没有那么多的语法束缚，用起来容易，学起来也很简单。

当要实际完成一项编程任务时，肯定需要借助各种工具，Python 提供了非常丰富的工具包来解决各种数据处理、分析、建模等问题。我们只要调用工具包，就可以轻轻松松地完成工作，相当于前人已经种好了树，我们去乘凉就好了。

那么，Python 在其他领域应用得怎么样呢？大家可能听过“Python 全栈开发”这个概念，所以 Python 相当于“万金油”，只要把它学好了，应用还是十分广泛的。

总结起来就是一句话：简洁、高效，用起来舒服！对于初学者来说，Python 是很友好的，可以说它是最简单易学的编程语言。

1.1.2 Python 该怎么学

很多零基础的读者的第一想法可能就是先去买一本非常厚的 Python 教材，然后慢慢地从入门到精通……其实我认为语言只是用来帮助解决问题的工具，不建议去找一本特别厚的书，来个半年学习计划，用最短的

时间学习最基础的、暂时够用的知识就可以了，越高级的语法用到的概率越小，先入手用起来，然后边做案例边学习才是高效的学习方法。

推荐大家先熟悉 Python 的基础部分，到图书馆随便找这方面的书，或者看看 Python 的在线课程都可以，有其他语言基础的同学学习 2 ~ 3 天就能用起来，第一次接触编程语言的人花一周的时间也会学得差不多了。

在后续的章节中本书还会涉及 Python 工具包的使用，其实这些工具的使用方法在其官方文档中都写得清清楚楚，并不需要全部背下来，只需要熟练操作即可，真正用到它的时候，还是要看看文档中每一个参数的具体含义。

1.2 人工智能的核心——机器学习

到底该如何学习人工智能呢？可以说，人工智能这个圈子太大了，各行各业都有涉及，可选择的方向也五花八门、各不相同，包括数据挖掘、计算机视觉、自然语言处理等各大领域。那么，是不是每个方向要学习的内容差别很大呢？不是的。其实最核心的就是机器学习，你要做的一切都离不开它，所以无论选择哪个领域，一定要把基础打牢。因此，第一个目标就是搞定机器学习的各大算法，并掌握其应用实践方法。

1.2.1 什么是机器学习

可能有些读者对机器学习还不是很熟悉，只不过因为最近这个词比较火才准备投身这个领域中。举一个小例子，我以前特别喜欢玩一款叫作《梦幻西游》的游戏。弃坑之后，游戏方的客服经理总给我打电话，说“迪哥能不能回来接着玩耍（充值）呀，帮派的小伙伴都十分想念你……”这时候我就想：他们为什么会给我打电话呢？这款游戏每天都有用户流失，不可能给每个用户都打电话吧，那么肯定是挑重点用户来沟通了。其后台肯定有玩家的各种数据，例如游戏时长、充值金额、战斗力等，通过这些数据就可以建立一个模型，用来预测哪些用户最有可能返回来接着玩啦！

机器学习要做的就是数据中学习有价值的信息，例如先给计算机一堆数据，告诉它这些玩家都是重点客户，让计算机去学习一下这些重点客户的特点，以便之后在海量数据中能快速将它们识别出来。

机器学习能做的远不止这些，数据分析、图像识别、数据挖掘、自然语言处理、语音识别等都是以其为基础的，也可以说人工智能的各种应用都需要机器学习来支撑（见图 1-2）。现在各大公司越来越注重数据的价值，人工成本也是越来越高，所以机器学习也就变得不可或缺了。

再给大家简单介绍一下学会机器学习之后可能从事的岗位，最常见的就是数据挖掘岗，即通过建立机器学习模型来解决实际业务问题，就业前景还是非常不错的，基本所有和数据打交道的公司都需要这个岗位。



图 1-2 机器学习的应用领域

接下来就是当下与人工智能结合最紧密的计算机视觉、自然语言处理和语音识别了。说白了就是要让计算机能看到、听到、读懂人类的数据。相对来说，我觉得计算机视觉领域发展会更快一些，因为随着深度学习技术的崛起，越来越多的研究人员加入这个行列，落地的项目更是与日俱增。自然语言处理和语音识别也是非常不错的方向，至于之后的路怎么走还是看大家的喜好吧，前提都是一样的——先把机器学习搞定！

1.2.2 机器学习的流程

上一小节简单介绍了机器学习的基本概念，那么机器学习是如何做事情的呢？下面通过一个简单例子来了解一下机器学习的流程（见图 1-3）。假设我们从网络上收集了很多新闻，有的是体育类新闻，有的是非体育类新闻，现在需要让机器准确地识别出新闻的类型。

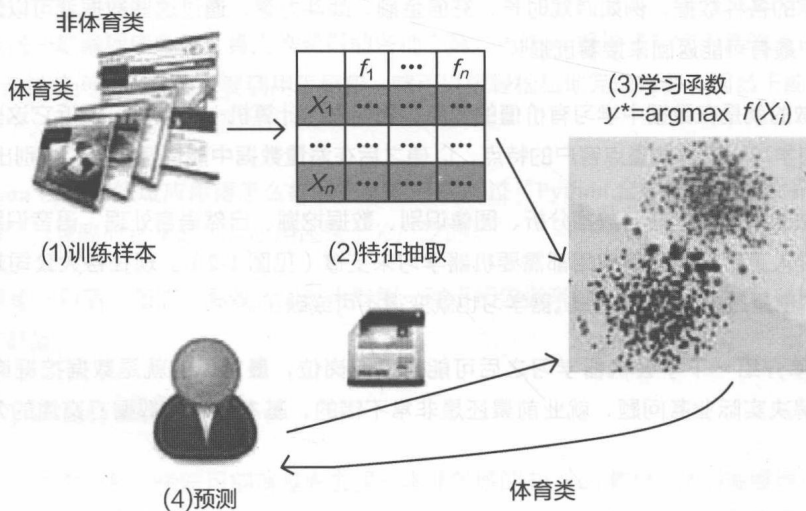


图 1-3 机器学习流程

一般来说，机器学习流程大致分为以下几步。

第①步：数据收集与预处理。例如，新闻中会掺杂很多特殊字符和广告等无关因素，要先把这些剔除掉。除此之外，可能还会用到对文章进行分词、提取关键词等操作，这些在后续案例中会进行详细分析。

第②步：特征工程，也叫作特征抽取。例如，有一段新闻，描述“科比职业生涯画上圆满句号，今天正式退役了”。显然这是一篇与体育相关的新闻，但是计算机可不认识科比，所以还需要将人能读懂的字符转换成计算机能识别的数值。这一步看起来容易，做起来就非常难了，如何构造合适的输入特征也是机器学习中非常重要的一部分。

第③步：模型构建。这一步只要训练一个分类器即可，当然，建模过程中还会涉及很多调参工作，随便建立一个差不多的模型很容易，但是想要将模型做得完美还需要大量的实验。

第④步：评估与预测。最后，模型构建完成就可以进行判断预测，一篇文章经过预处理再被传入模型中，机器就会告诉我们按照它所学数据得出的是什么结果。

1.2.3 机器学习该怎么学

很多读者可能都会有这种想法：工具包已经非常成熟了，是不是会调用工具包就可以了呢？笔者认为掌握算法原理与实际应用都是很重要的，很多人容易忽略算法的推导，这对之后的学习和应用肯定是不利的，因为做一件事情不能盲目去做，需要知道为什么要这么做！工具包也一样，不仅要学会使用它，更要知道其中每一个参数的作用，以及每一步操作在算法中都是什么含义。

这就需要熟悉每一个算法是怎么来的，每一步数学公式的目的是什么，数据是怎么一步步变成最后的决策结果的，每一步的参数又会对最终的结果产生什么样的影响。这几点都是非常重要的，所以在学习过程中需要深入其中每一步细节。

学习过程肯定有些枯燥，最好先从整体上理解其工作原理，然后再深入到每一处细节。这其中会涉及很多数学知识，对于初学者来说最头疼的就是这些公式和符号了，让大家从头到尾先学一遍数学可能有点不现实，所以遇到问题或者不理解的地方还需要大家勤动手，边学边查，也就是“哪里不会点哪里”。本书中所有知识点也都是按照笔者的理解跟大家分享的，所以不要惧怕数学，也不要过于钻牛角尖，理解即可。

1.3 环境配置

现在跟大家说一说本书所需的环境配置，也就是后续案例怎么玩起来，这个很重要，能给大家节省很多时间。我们要安装 Python 所需环境，不推荐去 Python 官网下载一个安装包，否则之后的配置和要安装的东西就太多了。

1.3.1 Anaconda 大礼包

环境配置时只需下载 Anaconda 即可，它相当于一个“全家桶”，里面不仅有 Python 所需环境，而且还把后续要用到的工具包和编程环境全部搞定了。首先登录 Anaconda 官网（<https://www.anaconda.com/download/>），下载对应软件，如图 1-4 所示。

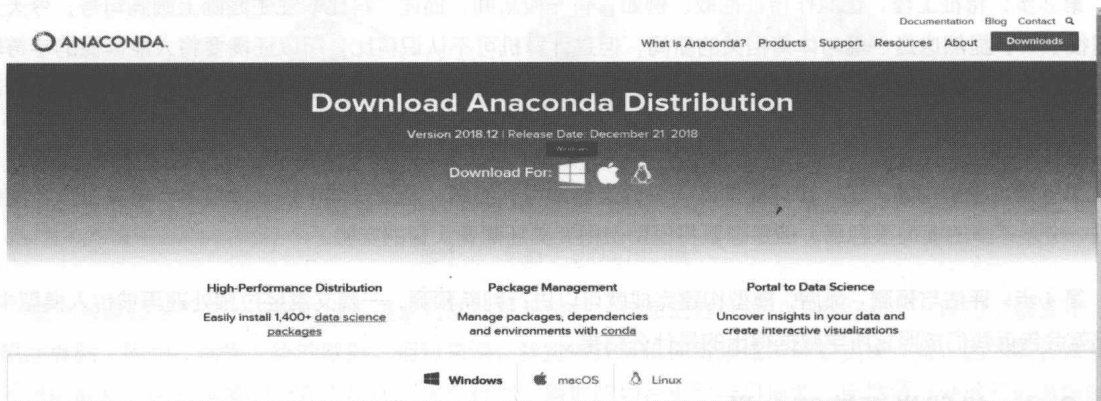


图 1-4 Anaconda 下载

然后根据自己的电脑选择不同的操作系统，并选择是 64 位的还是 32 位的。如果电脑是 32 位的，可以考虑换一换，因为很多工具包都不支持。

一定要选择 Python 3 版本（见图 1-5），几年前我在讲课和工作的時候用的是 Python 2.7 版本，当时，2.7 版本用的人比较多，而且相对稳定。但是从现在的角度出发，很多工具包都不支持 2.7 版本了，所以直接下载 3 版本即可。如果下载速度比较慢，读者可以登录镜像网址 <https://mirrors.tuna.tsinghua.edu.cn/anaconda/archive>，下载对应版本软件。

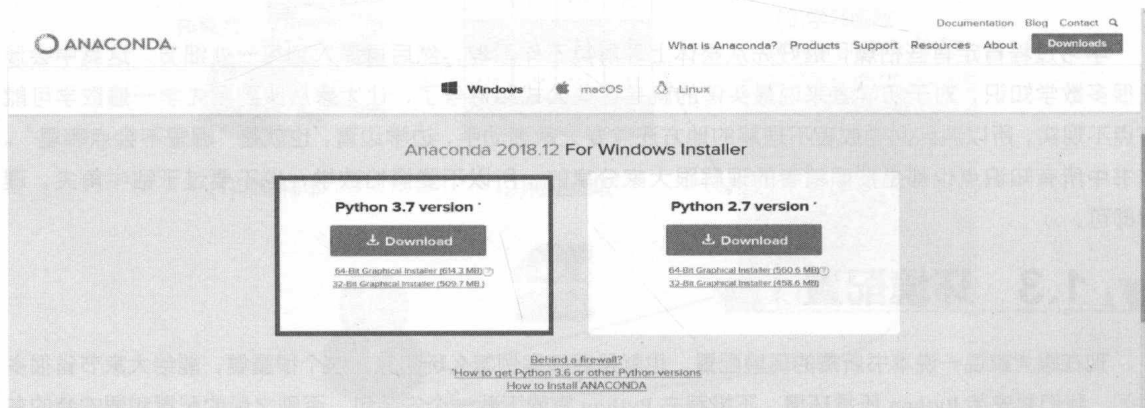


图 1-5 Python 版本选择