

ITエンジニアのための  
機械学習理論入門

# 机器学习 入门之道

[日]中井悦司◎著 姚待艳◎译



中国工信出版集团



人民邮电出版社  
POSTS & TELECOM PRESS

ITエンジニアのための  
機械学習理論入門

# 机器学习 入门之道

[日]中井悦司◎著 姚待艳◎译



人民邮电出版社  
北京

## 图书在版编目(CIP)数据

机器学习入门之道 / (日) 中井悦司著; 姚待艳译

北京: 人民邮电出版社, 2018. 5(2018. 9重印)

ISBN 978-7-115-47934-1

I. ①机… II. ①中… ②姚… III. ①机器学习

IV. ①TP181

中国版本图书馆CIP数据核字(2018)第032931号

### 内 容 提 要

人工智能正在形成一股新的浪潮,它将从技术、经济、社会等各个层面改变我们的工作和生活方式。作为实现人工智能的重要技术,机器学习正在受到人工智能专家之外的更广泛人群的关注,想要了解机器学习相关知识和技术的人日益增多。

本书紧紧围绕“机器学习的商业应用”这个主题,从数学原理上解释了机器学习的一些基础算法,如最小二乘法、最优推断法、感知器、Logistic回归、K均值算法、EM算法、贝叶斯推断等。全书的主旨在于帮助读者理解机器学习的本质,因此作者介绍具体的例题时,基本的着眼点是教会读者使用什么样的思维方式,以及如何计算,为读者探索更加复杂的深度学习领域或神经网络算法打下坚实的基础。

本书适合所有对机器学习感兴趣的读者阅读,尤其适合具备一定数学基础的IT工程师阅读,也可作为高等院校相关专业师生的参考读物。

---

◆ 著 [日] 中井悦司

译 姚待艳

责任编辑 陈宏

责任印制 焦志炜

◆ 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>

固安县铭成印刷有限公司印刷

◆ 开本: 700×1000 1/16

印张: 14

2018年5月第1版

字数: 180千字

2018年9月河北第2次印刷

著作权合同登记号 图字: 01-2016-5087

---

定价: 59.00元

读者服务热线: (010) 81055656 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

## 自序

---

差不多一年前，我就产生了这样一个怀疑——从事机器学习相关领域工作的 IT 工程师可能会以超过预想的速度不断增加。从“数据科学”到“深度学习”，甚至是“人工智能”，全都充斥着流行媒体，即使对不是从事数据分析的普通 IT 工程师来说，期待已久的机器学习应用时代也已经到来了。当今甚至还有机构宣称自己可以提供“无需专业知识也能使用”的机器学习服务。

但是，这里面有很大的陷阱。市面上各种关于机器学习的工具和程序库都是开源的，好像谁都可以掌握机器学习，输入数据、运行程序就可以得出一定的结果。可是，这样的结果到底有什么“意义”呢？既然是把机器学习的结果应用到商业领域，就必须理解其中的算法并正确掌握结果所传达的含义。

本书始终围绕“机器学习的商业应用”这个主题，从原理上解释了机器学习的基础算法。对于具体的例题，本书从“使用什么样的思维方式以及如何进行计算”的视角展开详细说明。因此，本书旨在帮助读者理解机器学习以及数据科学的本质。机器学习涵盖了各种各样的算法，本质上算法之间具有一个共通的思维方式，即“数据建模和参数调优”。本书把这种思维方式作为重点，尽量以通俗易懂的语言对各种数学公式进行说明。如果读者能够理解这种观点，那么面对本书没有涵盖的、更加复杂的深度学习或神经网络算法时，至少不会有畏难情绪。

“正在头疼于为委托方制定机器学习的商业应用方案”“突然决定参加销售额分析应用软件的开发项目”，我从 IT 工程师朋友那里听到了这样的声音，这引发了本文开头的那个怀疑。当今时代，如果能理解并熟练运用机器学习技术，IT 工程师们就能获得开启新的职业道路的机会。

最重要的是，机器学习中蕴含的趣味性可以满足 IT 工程师对知识的好奇心和对技术的探究之心。希望本书的读者都能以本书为起点，成功地踏出迈向机器学习世界的第一步。

2015 年初秋

中井悦司

## 写给读者的话

---

本书面向的主要读者是希望理解机器学习算法背后的原理并将机器学习灵活运用到商业领域中的 IT 工程师。机器学习有各种用途，本书将从数据分析的视角出发，对各种算法进行解释说明。但请各位读者注意，本书并非介绍机器学习工具和程序库使用方法的书籍。

本书采用的大部分例题都可以说是机器学习领域的经典例子，引用自下面的书籍：

C. M. Bishop. 模式识别和机器学习（上下册）. 东京：丸善出版社，2012.

虽然打开这本书的各位已经决定通过这本书学习机器学习，但很有可能一些读者并不能深入理解其中的原理。本书的读者当然并不局限于 IT 工程师，这本书也可以作为读者突破“权威经典”的入门书籍。

## 本书的阅读方法

本书将对机器学习的各种算法进行系统的介绍，请从第 1 章开始按顺序阅读。在第 1 章中，为了明确“机器学习的商业应用”这个主题，我会把机器学习放在体系更庞大的数据科学中进行介绍。在之后的第 2 章至第 8 章中，我将具体的算法应用到了第 1 章介绍过的具有代表性的例题中。我会对同一个问题应用不同的算法，这能让读者理解各种算法的特征并掌握通用的思考方法。

此外，本书还提供了用 Python 编写的各个算法的可运行示例代码。运行示例代码并观察具体的输出结果，就可以捕捉到数学公式无法呈现出来的算法的本质。

当然，要想理解机器学习的算法，就必须具备一定程度的数学知识。

本书尽量以“该数学公式用于什么计算”这种通俗易懂的方式进行解释，如果读者具备大学初级程度的数学知识，就能顺畅地理解书中内容。对机器学习相关的数学知识感兴趣的读者，可阅读微积分、线性代数和概率统计方面的书籍。

最后，针对“已经有将近 10 年没有碰过数学了”的读者，本书使用的主要数学符号和基本公式都归纳总结在了接下来的几页。需要的话，读者可以参考这部分内容。

## 致 谢

在此向在本书的撰写和出版过程中给我提供过帮助的各位表示衷心的感谢。

本书的策划萌芽于技术评论出版社池本公平先生的提案。感谢各位对我的机器学习相关知识体系进行整理，总结并出版了面向 IT 工程师的书籍，这是一个非常难得的机会。

我还要感谢本书的审校织学先生，他不仅快速高效地提出了修改意见，还提供了 Mac OS X/Windows 版本的安装步骤。

书中很多内容都受到了日本国立信息学研究所旗下项目“TopSE”的志愿者举办的学习讨论会的启发。在此感谢参与学习讨论会的各方人士。

编写本书时，很多时候我都是先把今年刚上小学的女儿送到车站，再到早间营业的星巴克进行写作。感谢照料我起居的妻子，她始终信奉“早睡早起身体好”，提倡健康的生活方式。我要对她说：“不胜感激！”

## 主要数学符号和基本公式

### ■ 求和符号

符号  $\Sigma$  表示求和。下式为  $x_1$  到  $x_N$  的求和运算:

$$\sum_{n=1}^N x_n = x_1 + x_2 + \cdots + x_N \quad (1)$$

### ■ 乘积符号

符号  $\Pi$  表示乘积。下式为  $x_1$  到  $x_N$  的乘法运算:

$$\prod_{n=1}^N x_n = x_1 \times x_2 \times \cdots \times x_N \quad (2)$$

### ■ 指数函数

符号  $\exp$  表示以自然常数  $e \approx 2.718$  为底的指数函数。下式表示  $e$  的  $x$  次方函数:

$$\exp x = e^x \quad (3)$$

指数函数的积可变换为自变量的和。

$$\prod_{n=1}^N e^{x_n} = e^{x_1} \times \cdots \times e^{x_N} = e^{x_1 + \cdots + x_N} = \exp \left\{ \sum_{n=1}^N x_n \right\} \quad (4)$$

指数函数  $e^x$  的微分形式函数不变。

$$\frac{d}{dx} e^x = e^x \quad (5)$$

## ■ 对数函数

符号  $\ln$  表示以自然常数  $e \approx 2.718$  为底的对数函数。

$$\ln x = \log_e x \quad (6)$$

代入  $x = e$ ，则其值变为 1。

$$\ln e = 1 \quad (7)$$

对数函数满足下面的对数法则：

$$\ln \frac{ab}{c} = \ln a + \ln b - \ln c \quad (8)$$

$$\ln a^b = b \ln a \quad (9)$$

因此，将公式 (4) 的指数函数代入对数函数中，可简化为下面的公式：

$$\ln \left( \exp \sum_{n=1}^N x_n \right) = \sum_{n=1}^N x_n \times \ln e = \sum_{n=1}^N x_n \quad (10)$$

该公式表示对数函数  $\ln x$  为指数函数  $e^x$  的反函数。

对数函数的微分形式如下：

$$\frac{d}{dx} \ln x = \frac{1}{x} \quad (11)$$

## ■ 偏微分

对于多变量函数，特定变量的微分称为偏微分（符号为  $\partial$ ）。

$$\frac{\partial f(x, y)}{\partial x} : y \text{ 不变, 对 } x \text{ 微分}$$

$\frac{\partial f(x, y)}{\partial y}$  :  $x$  不变, 对  $y$  微分

复合函数的微分公式对偏微分也成立。

$$\frac{\partial f(g(x, y))}{\partial x} = f'(g(x, y)) \times \frac{\partial g(x, y)}{\partial x} \quad (12)$$

$f'(x)$  表示一次微分系数。

$$f'(x) = \frac{df(x)}{dx} \quad (13)$$

## ■ 向量的内积和外积

公式中的粗体变量表示向量以及行列式。纵向排列的向量称为基础“列向量”。

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \quad (14)$$

根据书写方法的不同, 使用行向量进行描述时, 使用转置符号  $T$  表示列向量。

$$\mathbf{x} = (x_1, x_2, x_3)^T \quad (15)$$

反之, 对列向量进行转换可得到行向量。

$$\mathbf{x}^T = (x_1, x_2, x_3) \quad (16)$$

内积表示“行向量  $\times$  列向量”。

$$\mathbf{w}^T \mathbf{x} = (w_1, w_2, w_3) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \sum_{i=1}^3 w_i x_i \quad (17)$$

外积表示“列向量  $\times$  行向量”。

$$\mathbf{w}\mathbf{x}^T = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix} (x_1, x_2, x_3) = \begin{pmatrix} w_1 x_1 & w_1 x_2 & w_1 x_3 \\ w_2 x_1 & w_2 x_2 & w_2 x_3 \\ w_3 x_1 & w_3 x_2 & w_3 x_3 \end{pmatrix} \quad (18)$$

使用式(12),对代入向量内积的函数,可以求出特定元素的偏微分。

$$\frac{f(\mathbf{w}^T \mathbf{x})}{\partial w_i} = f'(\mathbf{w}^T \mathbf{x}) \frac{\partial(\mathbf{w}^T \mathbf{x})}{\partial w_i} = f'(\mathbf{w}^T \mathbf{x}) x_i \quad (19)$$

向量的大小用下面的符号表示:

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{x_1^2 + x_2^2 + x_3^2} \quad (20)$$

## ■ 随机变量的期望与方差

以一定概率取得各种不同的值的变量  $X$  称为随机变量,用  $P(x)$  表示  $X=x$  时的概率。随机变量的期望  $E$  和方差  $V$  由如下公式定义:

$$E[X] = \sum_x x P(x) \quad (21)$$

$$V[X] = E[\{X - E(X)\}^2] \quad (22)$$

式(21)中的和  $\sum_x$  为所有满足条件的  $x$  的和。

期望和方差之间有如下公式成立:

$$E[aX + b] = aE[X] + b \quad (23)$$

$$V[aX] = a^2 V[X] \quad (24)$$

$$V[X] = E[X^2] - (E[X])^2 \quad (25)$$

$\bar{x} = E[X]$ , 由式 (23) 可知下式成立:

$$E[X - \bar{x}] = E[X] - \bar{x} = 0 \quad (26)$$

假设有两个“独立的”随机变量  $X$  和  $Y$ ,  $X = x$  且  $Y = y$  时的概率 (同时发生的概率)  $P(x, y)$  可表示为  $X = x$  时的概率与  $Y = y$  时的概率的乘积。

$$P(x, y) = P_X(x) P_Y(y) \quad (27)$$

例如, 掷两个骰子同时掷出 1 点的概率, 为各个骰子分别掷出 1 点时的概率的乘积的  $1/6$ 。这表明不同的骰子掷出某个点的概率是相互独立的。

随机变量  $X$  和  $Y$  独立时,  $\bar{x} = E[X]$ ,  $\bar{y} = E[Y]$ , 则下式成立<sup>①</sup>:

$$\begin{aligned} E[(X - \bar{x})(Y - \bar{y})] &= \sum_{x, y} (x - \bar{x})(y - \bar{y}) P(x, y) \\ &= \sum_x (x - \bar{x}) P_X(x) \sum_y (y - \bar{y}) P_Y(y) \\ &= E[X - \bar{x}] E[Y - \bar{y}] = 0 \end{aligned} \quad (28)$$

---

① 这种关系会在 3.3.1 节中用到。

第 1 章

# 数据科学和机器学习

1

**1.1** 数据科学在商业领域中的作用 2

**1.2** 机器学习算法的分类 8

**1.2.1** 分类：产生类判定的算法 8

**1.2.2** 回归分析：预测数值的算法 9

**1.2.3** 聚类分析：对数据进行无监督群组化的算法 10

**1.2.4** 其他算法 12

**1.3** 本书使用的例题 13

**1.3.1** 基于回归分析的观测值推断 13

**1.3.2** 基于线性判别的新数据分类 17

**1.3.3** 图像文件的褪色处理（提取代表色） 18

**1.3.4** 识别手写文字 19

**1.4** 分析工具的准备 20

**1.4.1** 本书使用的数据分析工具 21

**1.4.2** 运行环境设置步骤（以 CentOS 6 为例） 22

**1.4.3** 运行环境设置步骤（以 Mac OS X 为例） 25

1.4.4 运行环境设置步骤（以 Windows 7/8.1 为例）—— 27

1.4.5 IPython 的使用方法—— 30

## 第 2 章

# 最小二乘法： 机器学习理论第一步 35

2.1 基于近似多项式和最小二乘法的推断—— 36

2.1.1 训练集的特征变量和目标变量—— 36

2.1.2 近似多项式和误差函数的设置—— 38

2.1.3 误差函数最小化条件—— 39

2.1.4 示例代码的确认—— 42

2.1.5 统计模型的最小二乘法—— 46

2.2 过度拟合检出—— 49

2.2.1 训练集和测试集—— 49

2.2.2 测试集的验证结果—— 50

2.2.3 基于交叉检查的泛化能力验证—— 52

2.2.4 基于数据的过度拟合变化—— 54

2.3 附录：Hessian 矩阵的特性—— 56

## 第 3 章

# 最优推断法： 使用概率的推断理论 59

3.1 概率模型的利用—— 60

3.1.1	“数据的产生概率”设置	60
3.1.2	基于似然函数的参数评价	65
3.1.3	示例代码的确认	69
<b>3.2</b>	<b>使用简化示例的解释说明</b>	<b>73</b>
3.2.1	正态分布的参数模型	74
3.2.2	示例代码的确认	76
3.2.3	推断量的评价方法（一致性和无偏性）	78
<b>3.3</b>	<b>附录：样本均值及样本方差一致性和无偏性的证明</b>	<b>80</b>
3.3.1	样本均值及样本方差一致性和无偏性的证明	81
3.3.2	示例代码的确认	85

## 第 4 章

### 感知器：

### 分类算法的基础

89

<b>4.1</b>	<b>概率梯度下降法的算法</b>	<b>91</b>
4.1.1	分割平面的直线方程	91
4.1.2	基于误差函数的分类结果评价	93
4.1.3	基于梯度的参数修正	95
4.1.4	示例代码的确认	99
<b>4.2</b>	<b>感知器的几何学解释</b>	<b>100</b>
4.2.1	对角项的任意性和算法的收敛速度	101
4.2.2	感知器的几何学解释	103
4.2.3	对角项的几何学意义	104

第 5 章

# Logistic 回归和 ROC 曲线： 学习模型的评价方法

107

<b>5.1</b>	对分类问题应用最优推断法	108
<b>5.1.1</b>	数据发生概率的设置	108
<b>5.1.2</b>	基于最优推断法的参数确定	112
<b>5.1.3</b>	示例代码的确认	114
<b>5.2</b>	基于 ROC 曲线的学习模型评价	117
<b>5.2.1</b>	Logistic 回归在实际问题中的应用	118
<b>5.2.2</b>	基于 ROC 曲线的性能评价	120
<b>5.2.3</b>	示例代码的确认	123
<b>5.3</b>	附录：IRLS 法的推导	126

第 6 章

# K 均值算法： 无监督学习模型的基础

133

<b>6.1</b>	基于 K 均值算法的聚类分析和应用实例	134
<b>6.1.1</b>	无监督学习模型类聚类分析	134
<b>6.1.2</b>	基于 K 均值算法的聚类分析	135
<b>6.1.3</b>	在图像数据方面的应用	138
<b>6.1.4</b>	示例代码的确认	141
<b>6.1.5</b>	K 均值算法的数学依据	143
<b>6.2</b>	“懒惰”学习模型 K 近邻法	146

6.2.1 基于  $K$  近邻法的分类 ————— 146

6.2.2  $K$  近邻法的问题 ————— 148

第 7 章

## EM 算法：

基于最优推断法的监督学习 151

7.1 使用伯努利分布的最优推断法 ————— 152

7.1.1 手写文字的合成方法 ————— 153

7.1.2 基于图像生成器的最优推断法应用 ————— 154

7.2 使用混合分布的最优推断法 ————— 157

7.2.1 基于混合分布的概率计算 ————— 157

7.2.2 EM 算法的过程 ————— 158

7.2.3 示例代码的确认 ————— 161

7.2.4 基于聚类分析的探索性数据解析 ————— 165

7.3 附录：手写文字数据的采集方法 ————— 167

第 8 章

## 贝叶斯推断：以数据为基础

提高置信度的手法 169

8.1 贝叶斯推断模型和贝叶斯定理 ————— 170

8.1.1 贝叶斯推断的思路 ————— 171

8.1.2 贝叶斯定理入门 ————— 172