

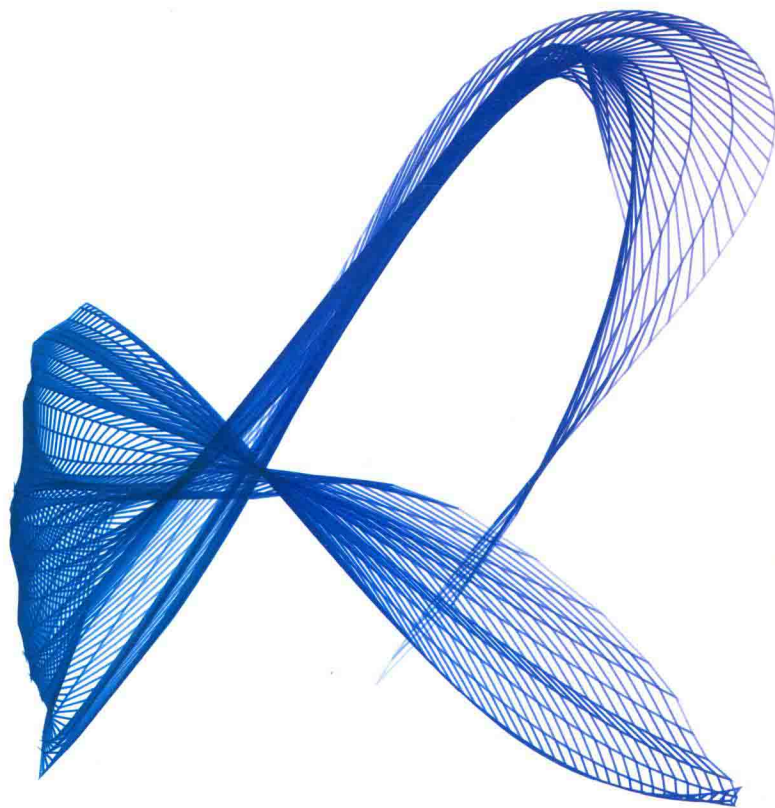


科大讯飞大数据专家团队撰写，不囿于Spark机器学习库，突出工程化思维与实践

6大算法模型构建，5大场景（异常检测、用户画像、点击率预估、企业征信、智慧交通）应用，从内涵认知到实践技能，全面提升



技术丛书



Spark机器学习 进阶实战

马海平 于俊 吕昕 向海◎著



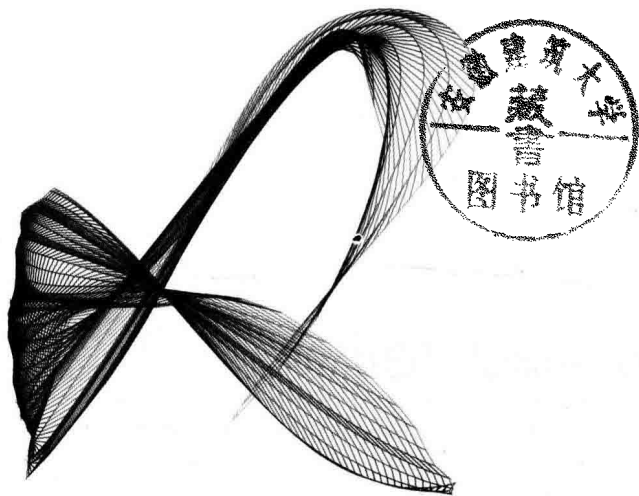
机械工业出版社
China Machine Press



技术丛书

Spark机器学习 进阶实战

马海平 于俊 吕昕 向海◎著



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

Spark 机器学习进阶实战 / 马海平等著. —北京: 机械工业出版社, 2018.9
(大数据技术丛书)

ISBN 978-7-111-60810-3

I. S… II. 马… III. 数据处理软件 IV. TP274

中国版本图书馆 CIP 数据核字 (2018) 第 201129 号

Spark 机器学习进阶实战

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: 高婧雅

责任校对: 李秋荣

印刷: 北京市兆成印刷有限责任公司

版次: 2018 年 9 月第 1 版第 1 次印刷

开本: 186mm × 240mm 1/16

印张: 14

书号: ISBN 978-7-111-60810-3

定价: 59.00 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



马海平

科大讯飞大数据研究院研究主管，中国科学与技术大学计算机技术博士，专注数据挖掘和人工智能算法的研究，及其在计算广告和个性化教育等方向的落地应用。



于俊

科大讯飞大数据专家，专注大数据和人工智能应用方案设计、基于Spark的大数据分析和价值挖掘，在大数据算法工程化实现方面具有丰富经验。



吕昕

科大讯飞大数据专家，专注大数据和人工智能技术在消费者业务中的应用、基于Spark的大数据分析和算法建模，在用户画像、内容推荐和精准营销领域有丰富的实践。



向海

邂智科技算法负责人，前科大讯飞大数据专家。专注Spark机器学习在智能客服中的应用，在NLP与对话机器人应用方面有丰富经验。

科大讯飞大数据专家团队撰写，不囿于Spark机器学习库，突出算法的工程化思维与实践。从基础引出算法，从算法实践到场景应用，层层推进，分享笔者的一些想法和见解，铺展开更为深入、全面的思路。

6大机器学习模型构建

分类：刻画事物特征的类标识，有效预测未知数据的归类情况。

聚类：根据相似程度生成对象集合，同集合相似，不同集合相异。

回归：找出数据规律和趋势，预测数据未来变化。

关联规则：挖掘关联关系，辅助商业决策。

协同过滤：刻画用户相似兴趣，实现偏好预测。

降维：有效地消除无关和冗余特征，提升模型精度。

5大典型应用场景

异常检测：有效解决入侵检测、欺诈检测、社交假新闻等问题。

用户画像：高度精炼用户的特征标识，为产品与决策提供数据支持和事实依据。

点击率预估：预估点击概率，计算点击收益，选出收益最高的策略。

企业征信：提供信用信息服务，洞察企业信用风险。

智慧交通：实现交通数据的价值，提供解决城市交通问题的思路。

同时，本书从《道德经》和《庄子》精选名言，并结合大数据机器学习相关内容，对名言加以讲解，引导大家以老庄的思想来认识大数据的内涵。

HZBOOKS | 华章IT | Information Technology



上善若水，水善利万物而不争。

数据一如水，无色无味，非方非圆，以百态存于自然，于自然无违也。绵绵密密，微则无声，巨则汹涌；与人无争却又容纳万物。生活离不开水，同样离不开数据，我们被数据包围，在数据中生活，体会着数据量爆炸式增长带来的幸福和挑战。

本书从《道德经》和《庄子》精选名言，并结合大数据机器学习相关内容，对名言加以讲解，引导大家以老庄的思想认识大数据的内涵，使用机器学习进行大数据价值挖掘，探求老子道之路和庄子智慧之路。

为什么要写这本书

2014年春天，曾经和公司大数据团队小伙伴一起聚焦研究大数据，为了解决国内资料匮乏、学习门槛较高的问题，着手编写《Spark 核心技术与高级应用》^①一书，并于2016年1月出版，取得了较好的反响，得到很多朋友的支持。

近年来，随着收集、存储和分析的数据量呈爆炸式增长，大规模的数据分析和数据价值挖掘能力已经成为影响企业生死存亡的关键，越来越多的企业必须面对这残酷而美好的挑战。基于大数据的机器学习有效解决了大数据带来的数据分析和数据挖掘瓶颈。

如何让更多的大数据从业人员更轻松地使用机器学习算法进行大数据价值挖掘，通过简单的学习建立大数据环境下的机器学习工程化思维，在不必深究算法细节的前提下，实现大数据分类、聚类、回归、协同过滤、关联规则、降维等算法，并使用这些算法解决实际业务场景的问题。2016年秋天，在机械工业出版社高婧雅编辑的指导下，怀着一颗附庸风雅之心，

① 该书已由机械工业出版社出版，书号为 ISBN 978-7-111-52354-3。——编辑注

我决定和小伙伴们一起朝着新的目标努力。

本书的写作过程中，Spark 版本也在不断变化，秉承大道至简的原则，我们一方面尽量按照新的版本进行统筹，另外一方面尽量做到和版本解耦，希望能抛砖引玉，以个人的一些想法和见解，为读者拓展出更深入、更全面的思路。

本书只是一个开始，如何使用机器学习算法从海量数据中挖掘出更多的价值，还需要无数的大数据从业人员前赴后继，突破漫漫雄关，共同创造美好的大数据机器学习时代。

本书特色

本书介绍大数据机器学习的算法和实践，同时对传统文化进行了一次缅怀，吸收传统文化的精华，精选了《道德经》和《庄子》部分名言，实现大数据和哲学思想的有效统一。结合老子的“无为”和庄子的“天人合一”思想，引导读者以辩证法思考方式认识大数据机器学习的内涵。

从技术层面上，本书一方面基于 Spark 现有的机器学习库讲解，另一方面尽量做到和现有 Spark 版本中的机器学习库解耦，突出对大数据机器学习的宏观理解，并给出典型算法的工程化实现，使更多的人轻松使用机器学习进行大数据价值挖掘，从而建立大数据机器学习工程化思维，在不必深究算法细节的前提下有效解决实际问题。本书更加强调在实际场景中的应用，并有针对性地给出了综合应用场景。

从适合读者阅读和掌握知识的结构安排上讲，本书分“基础篇”“算法篇”“综合应用篇”三个维度层层推进，便于读者在深入理解基础上根据相应的解决思路找到适合自己的方案。

本书使用的机器学习算法和应用场景都是实际业务的抽象，并基于具体业务进行实现。作为本书的延续，接下来我们会聚焦应用实践并提供更深层次的拓展，专注知识图谱的技术与应用，以及 Bot 技术与构建实战，期待相关图书能和读者尽早见面。

读者对象

(1) 对大数据感兴趣的读者

伴随着大数据时代的到来，很多工作都和大数据息息相关，无论是传统行业、互联网行业，还是移动互联网行业，都必须要了解大数据，通过大数据发现自身的价值。对这部分读者来说，本书的内容能够帮助他们加深对大数据 / 机器学习及其演进趋势的理解，通过本书可以了解机器学习相关算法，以及 Spark 机器学习应用场景和存在价值，如果希望更深层次地

掌握 Spark 机器学习相关知识，本书可以作为一个很好的开始。

(2) 从事大数据机器学习算法的研究人员

本书基于分类、聚类、回归、关联规则、协同过滤、降维等算法，结合异常检测、用户画像、广告点击率预估、企业征信大数据、智慧交通大数据等场景，系统地讲解了 Spark 机器学习相关知识，对从事大数据算法的研究人员来说，能够身临其境地体验各种场景，了解各类算法在不同场景下的优缺点，减少自己的研究成本。本书对生产环境中遇到的算法建模、数据挖掘等问题有很好的借鉴作用。

(3) 大数据工程开发人员

大数据工程开发人员可以从本书中获取需要的机器学习算法工程化知识。对大数据工程开发人员来说，掌握并快速对算法进行工程化，是很重要的技能，本书为填补算法工程开发人员与算法研究人员之间的鸿沟、高效工作提供了更多可能。

(4) 大数据架构设计人员

基于大数据的采集、存储、清洗、实时计算、统计分析、数据挖掘等是大数据架构师必备技能。他们需要对 Spark 机器学习进行了解，才能在架构设计中综合考虑各种因素，构建稳定高效的大数据架构。

如何阅读本书

本书分为三篇，共计 13 章内容。

基础篇（第 1 和 2 章），对机器学习进行概述讲解，并通过 Spark 机器学习进行数据分析。

算法篇（第 3 ~ 8 章），针对分类、聚类、回归、关联规则、协同过滤、降维等算法进行详细讲解，并进行算法建模应用实现。

综合应用篇（第 9 ~ 13 章），综合异常检测、用户画像，引出广告点击率预估，并对企业征信大数据、智慧交通大数据等场景进行实践，详细讲解基于 Spark 的大数据机器学习综合应用。

勘误和支持

由于笔者的水平有限，编写时间仓促，书中难免会出现一些错误或者不准确的地方，恳请读者批评指正。如果你有更多的宝贵意见，可以通过大数据技术交流 QQ 群 435263033，或者邮箱 datadance@163.com 联系我们，期待能够得到大家的真挚反馈，在大数据和人工智能征程中互勉共进。

致谢

感谢亲爱的搭档马海平、吕昕、向海三位大数据专家以及谭昶博士，在本书写作遇到困难的时候，我们一直互相鼓励，牺牲休息时间，坚持不放弃。

感谢大数据团队的张志勇、张龙、陈爱华、杨柳、俞祥祥、王庆庆、牛鑫、谢榭、李雅洁，以及廖攀、覃雪辉等小伙伴，你们为本书的修改贡献了宝贵的智慧，你们的参与使本书更上一层楼。

本书使用了部分互联网测试数据，包括：Stanford 的 gowalla 数据、360 的应用市场数据、UCI 的鸢尾花卉数据和裙子销售数据、数据堂的豆瓣电影评分数据、Digit 数据集、新闻 App 的用户行为数据、某运营商手机信令数据、某地图路况的道路拥堵指数数据，在这里进行特别感谢。

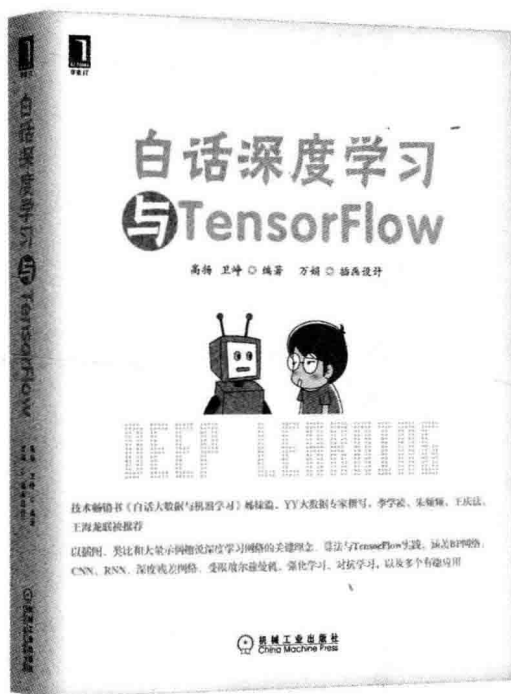
最后特别祝福本书写作期间出生的马海平家的二宝和向海家的二宝，你们的出生代表了大数据机器学习有了新的传承，也让我们的努力变得更有意义。

谨以此书献给大数据团队的小伙伴，以及众多热爱大数据机器学习技术的朋友！

于俊

2018年8月

推荐阅读

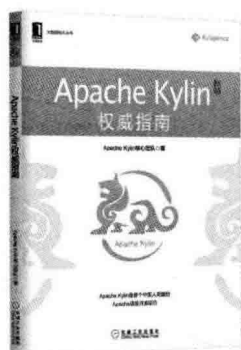
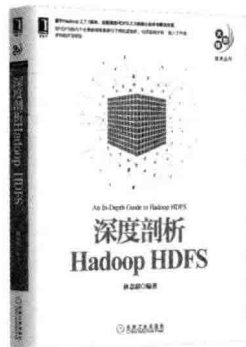


白话深度学习与TensorFlow

书号：978-7-111-57457-6 作者：高扬 定价：69.00元

《白话大数据与机器学习》姊妹篇，YY大数据专家多年实践沉淀之作，
简单易懂，轻松入门

推荐阅读



前 言

第一篇 基础篇

第 1 章 机器学习概述 2

1.1 机器学习概述 2

1.1.1 理解大数据 2

1.1.2 机器学习发展过程 4

1.1.3 大数据生态环境 5

1.2 机器学习算法 6

1.2.1 传统机器学习 6

1.2.2 深度学习 8

1.2.3 其他机器学习 8

1.3 机器学习分类 9

1.3.1 监督学习 9

1.3.2 无监督学习 10

1.3.3 半监督学习 10

1.3.4 强化学习 10

1.4 机器学习综合应用 11

1.4.1 异常检测 12

1.4.2 用户画像 12

1.4.3 广告点击率预估 12

1.4.4 企业征信大数据应用 12

1.4.5 智慧交通大数据应用 13

1.5 本章小结 13

第 2 章 数据分析流程和方法 14

2.1 数据分析概述 14

2.2 数据分析流程 15

2.2.1 业务调研 16

2.2.2 明确目标 16

2.2.3 数据准备 16

2.2.4 特征处理 17

2.2.5 模型训练与评估 21

2.2.6 输出结论 23

2.3 数据分析的基本方法 24

2.3.1 汇总统计 24

2.3.2 相关性分析 25

2.3.3 分层抽样 26

2.3.4 假设检验 26

2.4 简单的数据分析实践 27

2.4.1 环境准备 27

2.4.2 准备数据 28

2.4.3 数据分析 29

2.5 本章小结	30	4.2.1 KMeans 聚类	54
第二篇 算法篇		4.2.2 DBSCAN 聚类	55
第3章 构建分类模型		4.2.3 主题聚类	56
3.1 分类模型概述	32	4.3 聚类效果评价	58
3.2 分类模型算法	34	4.3.1 集中平方误差和	58
3.2.1 逻辑回归	34	4.3.2 Purity 评价法	59
3.2.2 朴素贝叶斯模型	36	4.4 使用 KMeans 对鸢尾花卉数据集聚类	59
3.2.3 SVM 模型	37	4.4.1 准备数据	59
3.2.4 决策树模型	39	4.4.2 特征处理	60
3.2.5 K-近邻	40	4.4.3 聚类分析	60
3.3 分类效果评估	40	4.4.4 模型性能评估	62
3.3.1 正确率	41	4.5 使用 DBSCAN 对 GPS 数据进行聚类	62
3.3.2 准确率、召回率和 F1 值	41	4.5.1 准备数据	63
3.3.3 ROC 和 AUC	42	4.5.2 特征处理	64
3.4 App 数据的分类实现	44	4.5.3 聚类分析	64
3.4.1 选择分类器	44	4.5.4 模型参数调优	65
3.4.2 准备数据	45	4.6 其他模型	66
3.4.3 训练模型	46	4.6.1 层次聚类	66
3.4.4 模型性能评估	48	4.6.2 基于图的聚类	67
3.4.5 模型参数调优	49	4.6.3 混合聚类模型	67
3.5 其他分类模型	50	4.7 本章小结	68
3.5.1 随机森林	50	第5章 构建回归模型	
3.5.2 梯度提升树	51	5.1 常用回归模型	69
3.5.3 因式分解机模型	51	5.1.1 线性回归模型	70
3.6 本章小结	52	5.1.2 回归树模型	70
第4章 构建聚类模型		5.1.3 其他回归模型	71
4.1 聚类概述	53	5.2 评估指标	73
4.2 聚类模型	54	5.3 回归模型优化	74

5.3.1	特征选择	74	第 7 章 协同过滤	97	
5.3.2	特征变换	74	7.1	协同过滤概述	97
5.4	构建 UCI 裙子销售 数据回归模型	75	7.2	常用的协同过滤算法	98
5.4.1	准备数据	75	7.2.1	基于用户的协同过滤	99
5.4.2	训练模型	78	7.2.2	基于物品的协同过滤	100
5.4.3	评估效果	79	7.2.3	矩阵分解技术	101
5.4.4	模型优化	79	7.2.4	推荐算法的选择	102
5.5	其他回归模型案例	80	7.3	评估标准	103
5.5.1	GDP 影响因素分析	81	7.3.1	准确率	103
5.5.2	大气污染分析	81	7.3.2	覆盖率	103
5.5.3	大数据比赛中的回归问题	81	7.3.3	多样性	104
5.6	本章小结	82	7.3.4	其他指标	104
第 6 章 构建关联规则模型		83	7.4	使用电影评分数据进行 协同过滤实践	104
6.1	关联规则概述	83	7.4.1	准备数据	105
6.2	常用关联规则算法	84	7.4.2	训练模型	106
6.2.1	Apriori 算法	84	7.4.3	测试模型	109
6.2.2	FP-Growth 算法	85	7.4.4	使用 ALS 结果	111
6.3	效果评估和优化	86	7.5	本章小结	112
6.3.1	效果评估	86	第 8 章 数据降维	113	
6.3.2	效果优化	87	8.1	降维概述	113
6.4	使用 FP-Growth 对豆瓣 评分数据进行挖掘	88	8.2	常用降维算法	114
6.4.1	准备数据	89	8.2.1	主成分分析	114
6.4.2	训练模型	89	8.2.2	奇异值分解	116
6.4.3	观察规则	91	8.2.3	广义降维	117
6.4.4	参数调优	91	8.2.4	文本降维	118
6.4.5	使用算法	92	8.3	降维评估标准	121
6.5	其他应用场景	94	8.4	使用 PCA 对 Digits 数据集 进行降维	122
6.6	本章小结	96	8.4.1	准备数据	122

8.4.2	训练模型	123
8.4.3	分析降维结果	124
8.5	其他降维方法	124
8.5.1	线性判别分析	124
8.5.2	局部线性嵌入	125
8.5.3	拉普拉斯特征映射	125
8.6	本章小结	126

第三篇 综合应用篇

第9章	异常检测	128
9.1	异常概述	128
9.1.1	异常的产生	129
9.1.2	异常检测的分类	129
9.2	异常检测方法	130
9.2.1	基于模型的方法	130
9.2.2	基于邻近度的方法	131
9.2.3	基于密度的方法	132
9.2.4	基于聚类的方法	133
9.3	异常检测系统	133
9.3.1	异常检测过程	133
9.3.2	异常检测步骤	134
9.3.3	特征选取和设计	135
9.4	应用场景	137
9.4.1	入侵检测	137
9.4.2	欺诈检测	138
9.4.3	社交假新闻	140
9.4.4	医疗和公共卫生	141
9.5	新闻 App 数据异常检测实践	141
9.5.1	准备数据	141
9.5.2	数据预处理	142

9.5.3	异常检测	142
9.6	本章小结	144

第10章 用户画像

10.1	用户画像概述	145
10.1.1	什么是用户画像	145
10.1.2	为什么需要用户画像	146
10.2	用户画像流程	147
10.2.1	整体流程	147
10.2.2	标签体系	148
10.3	构建用户画像	150
10.3.1	人口属性画像	150
10.3.2	兴趣画像	152
10.3.3	地理位置画像	155
10.4	用户画像评估和使用	155
10.4.1	效果评估	156
10.4.2	用户画像使用	157
10.5	新闻 App 用户画像实践	158
10.5.1	事实标签构建	158
10.5.2	兴趣标签构建	159
10.6	本章小结	161

第11章 广告点击率预估

11.1	点击率预估概述	162
11.1.1	互联网广告的发展	163
11.1.2	互联网广告交易架构	163
11.1.3	点击率预估应用	165
11.2	点击率预估技术	166
11.2.1	数据收集	166
11.2.2	特征构建	167
11.2.3	特征处理和选择	169

11.2.4	模型训练	170	12.3.1	企业信用报告	186
11.3	模型效果评估	172	12.3.2	企业风控管理	187
11.3.1	模型指标评估	172	12.4	企业法人资产建模实践	188
11.3.2	线上流量评估	172	12.4.1	建模流程	188
11.4	新闻 App 点击率预估实践	173	12.4.2	数据准备	190
11.4.1	特征提取	173	12.4.3	模型工程实现	191
11.4.2	模型训练	174	12.5	本章小结	194
11.4.3	广告 CTR 模型扩展	175			
11.5	本章小结	177			
第 12 章 企业征信大数据应用			第 13 章 智慧交通大数据应用		
12.1	征信概述	178	13.1	智慧交通大数据概述	195
12.1.1	征信组成	179	13.2	人群生活模式划分	196
12.1.2	传统征信	180	13.2.1	数据介绍	196
12.1.3	大数据征信	180	13.2.2	数据预处理	196
12.2	企业征信大数据平台	181	13.2.3	特征构建	197
12.2.1	大数据征信平台架构	181	13.2.4	生活模式挖掘	200
12.2.2	企业征信服务流程	182	13.2.5	划分结果分析	202
12.2.3	企业征信数据源	182	13.3	道路拥堵模式聚类	204
12.2.4	企业征信画像库	183	13.3.1	数据介绍	204
12.2.5	征信评分模型	185	13.3.2	数据预处理	205
12.3	企业征信大数据应用	186	13.3.3	特征构建	206
			13.3.4	拥堵模式挖掘	207
			13.4	本章小结	210