



德阳科贸职业学院
Deyang Vocational College of Technology and Trade

德阳科贸职业学院

《Python 数据采集与分析 专周实训》指导书

二级学院（部）： 信息工程学院

适 用 对 象： 2021 级大数据技术专业

执 笔 人： 唐雪

审 核 人： 张玮

制 定 时 间： 2023 年 2 月

德阳科贸职业学院教务处制

二〇二三年 二 月



目录

电影数据的处理与可视化	1
目标	1
软硬件环境	1
实训任务	1
实训步骤	5
参考代码	错误！未定义书签。

电影数据的处理与可视化

目标

1. 掌握 Python 的语法规则和基本的程序设计方法，能够使用 Python 实现对文档的读取、合并功能；
2. 掌握使用 Python 操作文档、处理数据的方法，能够开发简单的逻辑功能模块，实现对数据的处理功能；
3. 掌握数据可视化开发的方法，能够开发简单数据可视化功能。

软硬件环境

硬件：PC 电脑一台

配置：win7 或 win10 系统，内存大于 4G，硬盘 250G 及以上；

软件环境：Pycharm、Python3.7。

实训任务

本案例参考参考 Gred Reda 的分析

<http://grouplens.org/datasets/movielens/>

对电影评分数据集进行数据处理与可视化，设计要求如下：

一、需求分析

下载公开数据集《电影数据》，进行数据分析。

1. 理解数据，数据如下：



u.user

文档中数据如下:

	A	B	C	D	
1	196	242	3	881250949	1 24 M technician 85711
2	186	302	3	891717742	2 53 F other 94043
3	22	377	1	878887116	3 23 M writer 32067
4	244	51	2	880606923	4 24 M technician 43537
5	166	346	1	886397596	5 33 F other 15213
6	298	474	4	884182806	6 42 M executive 98101
7	115	265	2	881171488	7 57 M administrator 91344
8	253	465	5	891628467	8 36 M administrator 05201
9	305	451	3	886324817	9 29 M student 01002
10	6	86	3	883603013	10 53 M lawyer 90703
11	62	257	2	879372434	11 39 F other 30329
12	286	1014	5	879781125	12 28 F other 06405
13	200	222	5	876042340	13 47 M educator 29206
14	210	40	3	891035994	14 45 M scientist 55106
15	224	29	3	888104457	15 49 F educator 97301
16	303	785	3	879485318	16 21 M entertainment 10309
17	122	387	5	879270459	17 30 M programmer 06355
18	194	274	2	879539794	18 35 F other 37212
19	291	1042	4	874834944	19 40 M librarian 02138
20	234	1184	2	892079237	20 42 F homemaker 95660
21	119	392	4	886176814	21 26 M writer 30068
22	167	486	4	892738452	22 25 M writer 40206
23	299	144	4	877881320	23 30 F artist 48197
24	291	118	2	874833878	24 21 F artist 94533
25	308	1	4	887736532	25 39 M engineer 55107
26	95	546	2	879196566	

1	Toy Story (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Toy%20story%20(1995)]	[o][o][o][i][l][i][o][o][o][o][o][o][o][o][o][o][o][o]
2	GoldenEye (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?GoldenEye%20(1995)]	[o][i][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
3	Four Rooms (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
4	Get Shorty (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?GetShorty%20(1995)]	[o][i][o][o][o][i][o][o][o][o][o][o][o][o][o][o]
5	Copcat (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Copcat%20(1995)]	[o][o][o][o][o][i][o][o][o][o][o][o][o][o][o][o]
6	Shanghai Triad (Yao a yao dao wai po qiao) (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Yao+a+yao+dao+wai+po+qiao+(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
7	Twelve Monkeys (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
8	Babe (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Babe%20(1995)]	[o][o][o][o][o][i][o][o][o][o][o][o][o][o][o][o]
9	Dead Man Walking (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?DeadManWalking%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
10	Richard III (1995)	[22-Jan-1996]	[http://us.imdb.com/M/title-exact?RichardIII%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
11	Seven (Se7en) (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Se7en%20(1995)]	[o][o][o][o][o][i][o][o][o][o][o][o][o][o][o][o]
12	Usual Suspects, The (1995)	[14-Aug-1995]	[http://us.imdb.com/M/title-exact?UsualSuspects,%20The%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
13	Nighty Aphrodite (1995)	[30-Oct-1995]	[http://us.imdb.com/M/title-exact?NightyAphrodite%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
14	Postino, Il (1995)	[01-Jan-1994]	[http://us.imdb.com/M/title-exact?Postino,%20Il%20(1994)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
15	Fred, Holland & Son (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Fred,Holland&Son%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
16	Frank Trist (Garon Audit) (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?FrankTrist(GaronAudit)%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
17	Freaky Frank Till Dawn (1996)	[05-Feb-1996]	[http://us.imdb.com/M/title-exact?FreakyFrankTillDawn%20(1996)]	[o][i][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
18	White Balloon, The (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?WhiteBalloon,%20The%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
19	Antonia's Line (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?Antonia%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
20	Angels and Insects (1995)	[01-Jan-1995]	[http://us.imdb.com/M/title-exact?AngelsandInsects%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]
21	Muppet Treasure Island (1996)	[16-Feb-1996]	[http://us.imdb.com/M/title-exact?MuppetTreasureIsland%20(1996)]	[o][i][i][o][o][o][o][o][o][o][o][o][o][o][o][o]
22	Braveheart (1995)	[16-Feb-1996]	[http://us.imdb.com/M/title-exact?Braveheart%20(1995)]	[o][o][o][o][o][o][o][o][o][o][o][o][o][o][o][o]

2.将数据表读取并设置表头:



	user_id	age	sex	occupation	zip_code
0	1	24	M	technician	85711
1	2	53	F	other	94043
2	3	23	M	writer	32067
3	4	24	M	technician	43537
4	5	33	F	other	15213

	user_id	movie_id	rating	unix-timestamp
0	196	242	3	881250949
1	186	302	3	891717742
2	22	377	1	878887116
3	244	51	2	880606923
4	166	346	1	886397596

	movie_id	...	unix-imdb_url
0	1	...	http://us.imdb.com/M/title-exact?Toy%20Story%2...
1	2	...	http://us.imdb.com/M/title-exact?GoldenEye%20(...
2	3	...	http://us.imdb.com/M/title-exact?Four%20Rooms%...
3	4	...	http://us.imdb.com/M/title-exact?Get%20Shorty%...
4	5	...	http://us.imdb.com/M/title-exact?Copycat%20(1995)

3.合并上述三张数据表

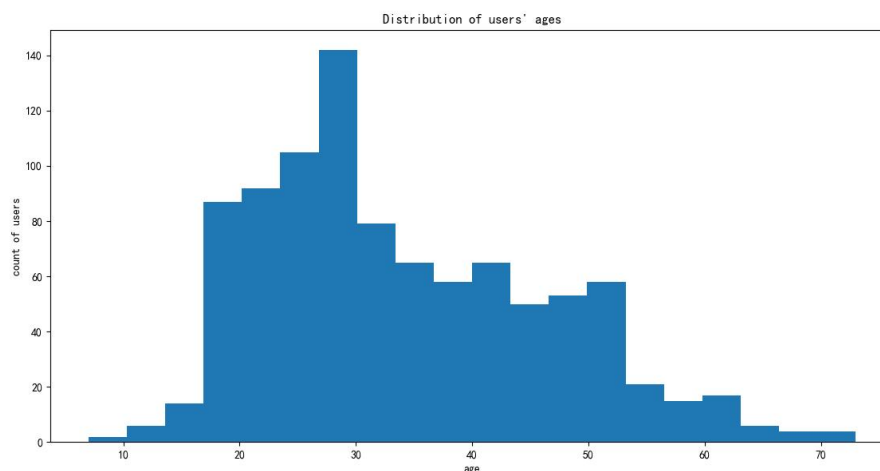
	movie_id	title	release_date	...	sex	occupation	zip_code
0	1	Toy Story (1995)	01-Jan-1995	...	M	retired	95076
1	4	Get Shorty (1995)	01-Jan-1995	...	M	retired	95076
2	5	Copycat (1995)	01-Jan-1995	...	M	retired	95076
3	7	Twelve Monkeys (1995)	01-Jan-1995	...	M	retired	95076
4	8	Babe (1995)	01-Jan-1995	...	M	retired	95076

4.按照四个数据分析指标进行电影数据评分的分析任务

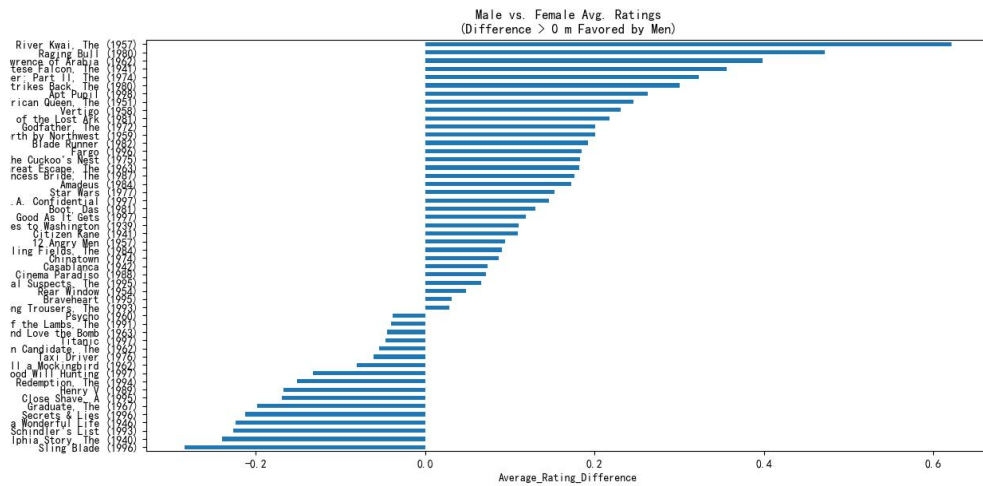


title	
Star Wars (1977)	583
Contact (1997)	509
Fargo (1996)	508
Return of the Jedi (1983)	507
Liar Liar (1997)	485
English Patient, The (1996)	481
Scream (1996)	478
Toy Story (1995)	452
Air Force One (1997)	431
Independence Day (ID4) (1996)	429
Raiders of the Lost Ark (1981)	420
Godfather, The (1972)	413
Pulp Fiction (1994)	394
Twelve Monkeys (1995)	392
Silence of the Lambs, The (1991)	390
Jerry Maguire (1996)	384
Chasing Amy (1997)	379

5.根据上一步中统计的数据，开发可视化功能，使用柱状图的方式展示评分与年龄的关系，如下图：



6.根据上一步中统计的数据，开发可视化功能，使用柱状图的方式展示评分与性别的关系，如下图：



二、实训步骤

1. 分析功能需求，设计功能模块；
2. 创建工程，安装工程使用的库；
3. 分析理解数据和业务指标；
4. 将数据进行预处理和特征提取；
5. 将得到的数据进行可视化展示