

## Python 期末大作业：分析新冠疫情数据

### 目录

一. 题目要求.....	2
二. 数据来源.....	2
三. 数据分析和展示.....	3
1. 15 天中，全球新冠疫情的总体变化趋势.....	3
2. 累计确诊数排名前 20 的国家名称及其数量.....	3
3. 15 天中，每日新增确诊数累计排名前 10 个国家的每日新增确诊数据的曲线图.....	4
4. 累计确诊人数占国家总人口比例最高的 10 个国家.....	4
5. 死亡率（累计死亡人数/累计确诊人数）最低的 10 个国家.....	5
6. 用饼图展示各个国家的累计确诊人数的比例.....	5
7. 展示全球各个国家累计确诊人数的箱型图，要有平均值.....	6
8. 其他展示.....	6
四. 根据以上数据，列出全世界应对新冠疫情最好的 10 个国家，并说明你的理由.....	7
五. 针对全球累计确诊数，利用前 10 天采集到的数据做后 5 天的预测.....	8
六. 核心代码.....	10
1. 爬虫程序（使用 selenium）.....	10
2. 各国应对措施打分代码.....	14
3. 数据展示代码.....	15
4. 数据预测代码.....	25

# 一. 题目要求

分析新冠疫情数据：找一个有全球新冠病毒数据的网站，爬取其中的数据（禁止使用数据接口直接获取数据）。要求爬取从 2020 年 12 月 1 日开始的连续 15 天的数据，国家数不少于 100 个。

# 二. 数据来源

初步观察发现，各网站的疫情数据各不相同，为保证最终数据分析不受网站因素影响，爬取数据使用的数据来源有三个。并在爬取到一定量的数据后确定使用哪个数据源。

一是百度疫情实时大数据报告(<https://voice.baidu.com/act/newpneumonia/newpneumonia>),一个是全球新冠疫情大数据分析平台(<https://www.zq-ai.com/#/fe/xgfybigdata>)，最后一个 是腾讯(<https://news.qq.com/zt2020/page/feiyang.htm#/global?nojump=1>)。首页截图如图 1：

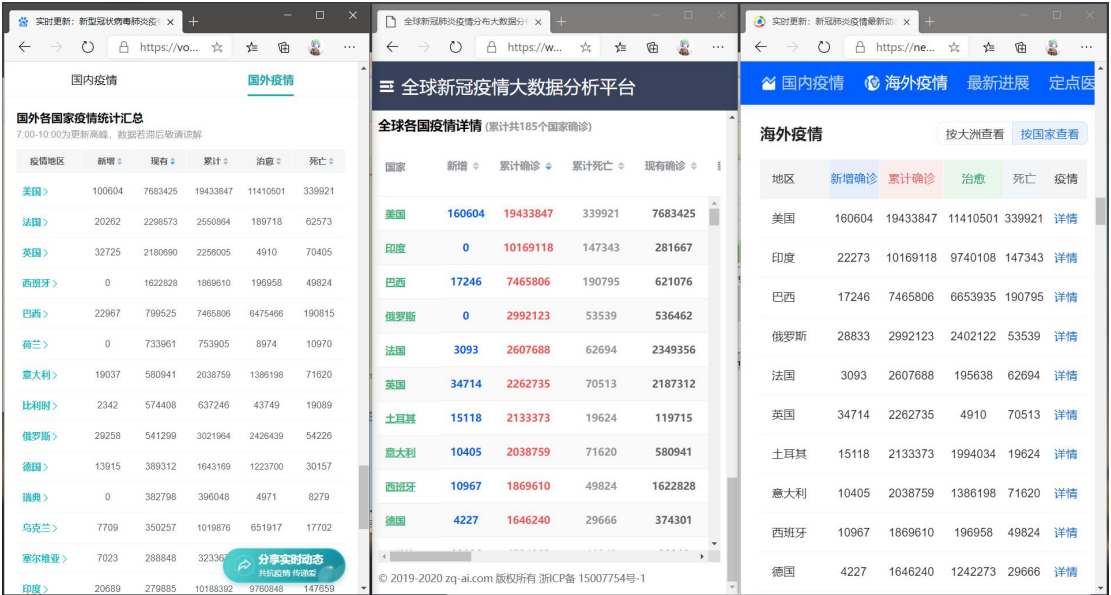


图 1 数据来源：左图是百度，中图是大数据分析平台，右图是腾讯

进行数据分析时，首先根据数据质量判断要使用哪个数据源。通过作图可以看到，百度的数据更新存在明显的滞后以及更新不及时的情况，就如下图图 2 中前三天数据无任何变化。与此同时也看到另一个数据源无此类问题。因此决定使用第二个网站的数据。



图 2 百度数据出现更新停滞的现象

此外，值得说明的是，为保证数据完全更新，每天爬取数据的时间为晚上 11 点 45 分。

### 三. 数据分析和展示

#### 1. 15 天中，全球新冠疫情的总体变化趋势



图 3 15 天中，全球新冠疫情的总体变化趋势

通过以上两幅线图可以直观地认识到，全球累计确诊数的总体变化趋势基本符合线性增长；每日新增在 50000 与 70000 之间波动。

#### 2. 累计确诊数排名前 20 的国家名称及其数量

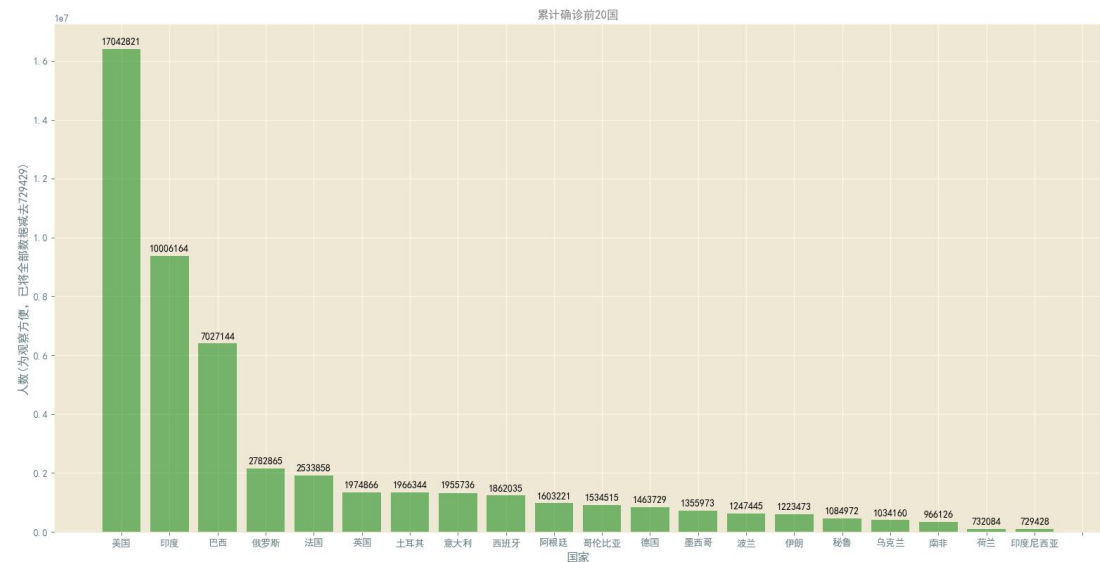


图 4 累计确诊数排名前 20 的国家名称及其数量

由上图的柱状图看到，累计确诊最多的是美国，其次是印度、巴西，这三个国家占据了累计确诊数的大部分；接着是一些欧洲国家如俄罗斯、法国、英国。

3. 15 天中，每日新增确诊数累计排名前 10 个国家的每日新增确诊数据的曲线图

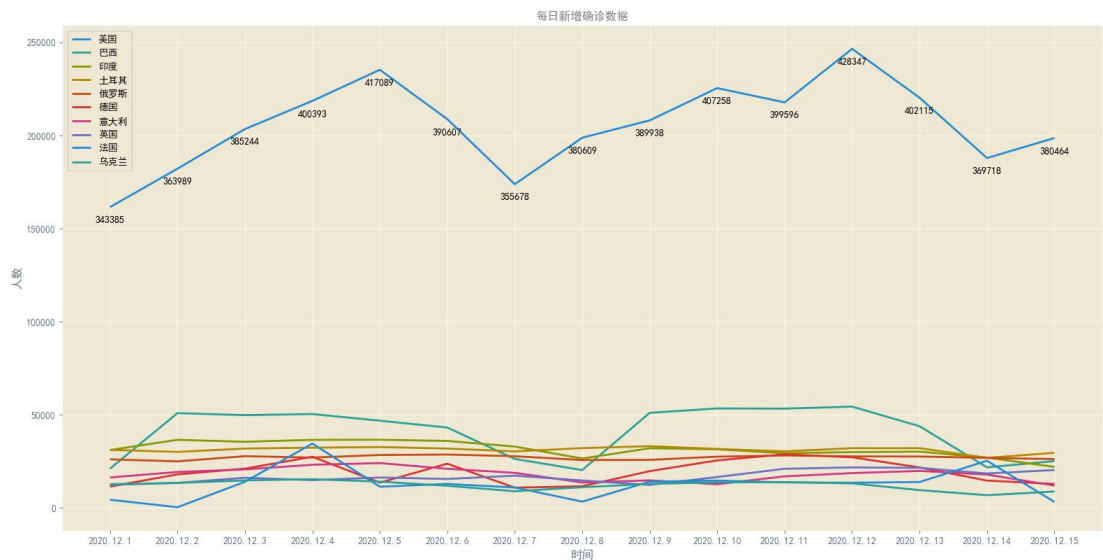


图 5 15 天中，每日新增确诊数累计排名前 10 个国家的每日新增确诊数据的曲线图

由于国家过于密集，故只在距离其他国家较远的“美国”曲线上标注出具体数值。

在曲线图中也可以看出，美国的每日新增远大于其他国家，对照图 3 的全球每日新增数据发现每天的新增确诊人数有约 60%来自美国。

此外，观察图 3 与图 5 的曲线，其大致走向是相同的。

4. 累计确诊人数占国家总人口比例最高的 10 个国家

通过数据处理，得到占比最高的国家是：安道尔，卢森堡，黑山，圣马力诺，巴林，比利时，卡塔尔，美国，亚美尼亚，格鲁吉亚。

并对这 10 个国家做出柱状图：

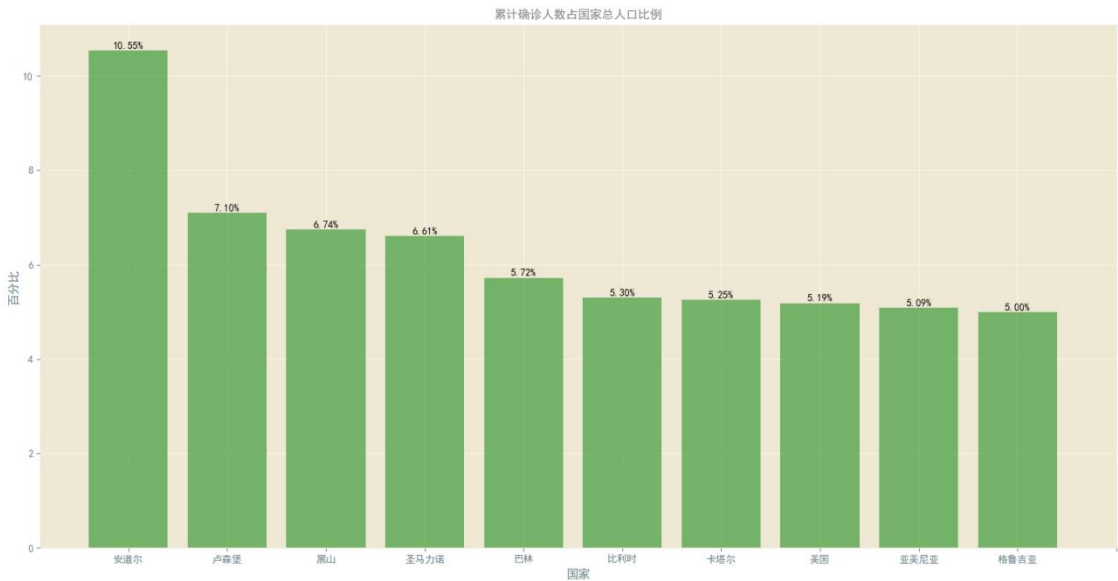


图 6 累计确诊人数占国家总人口比例最高的 10 个国家

以第一个国家安道尔为例，猜测其确诊人口占比高的重要因素是国土面积小以及人口数

少，导致一旦发生病例就会迅速波及到其大量人口。

5. 死亡率（累计死亡人数/累计确诊人数）最低的 10 个国家

通过数据处理，得到占比最高的国家是：新加坡，布隆迪，卡塔尔，博茨瓦纳，阿联酋，马尔代夫，巴林，摩纳哥，斯里兰卡，马来西亚。

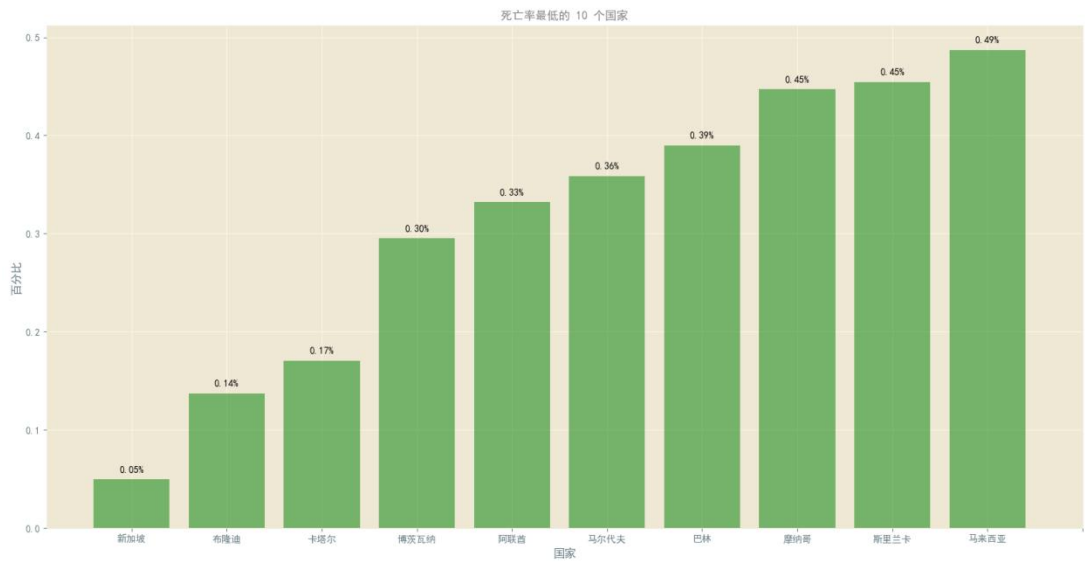


图 7 死亡率（累计死亡人数/累计确诊人数）最低的 10 个国家

6. 用饼图展示各个国家的累计确诊人数的比例

各个国家的累计确诊人数的比例 (仅展示前10个国家，剩下的归于“其他”)

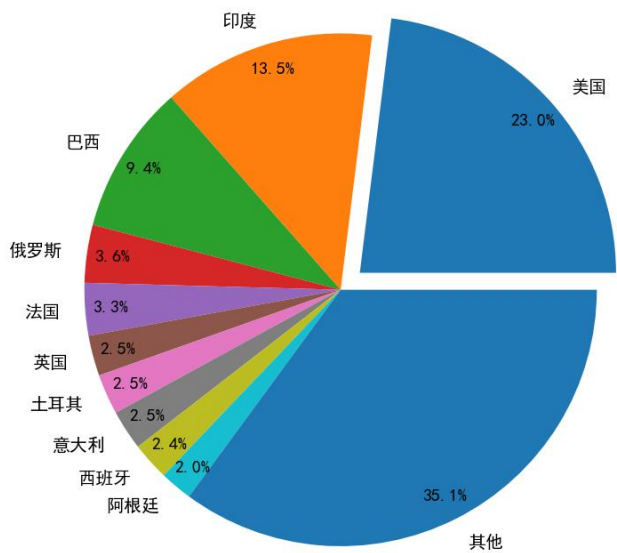


图 8 各个国家的累计确诊人数比例的饼图

由饼图可以更加直观地看出，美国、印度、巴西、俄罗斯四国的累计确诊合计达到全球的一半。

7. 展示全球各个国家累计确诊人数的箱型图，要有平均值

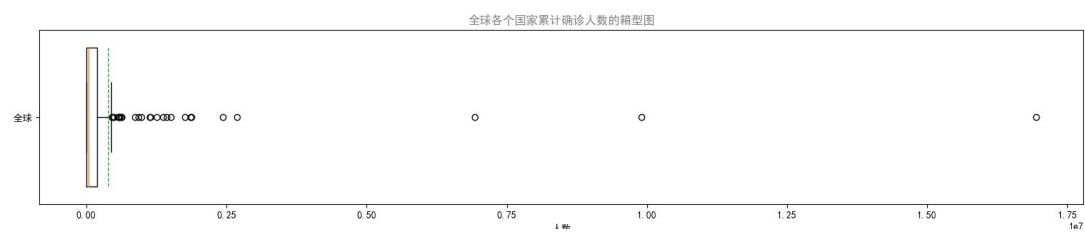


图 9 含异常值的箱型图

由于各国累计确诊差距较大，故异常值较多，影响观感，将上图去掉异常值并添加必要的注释后，绘图如下：

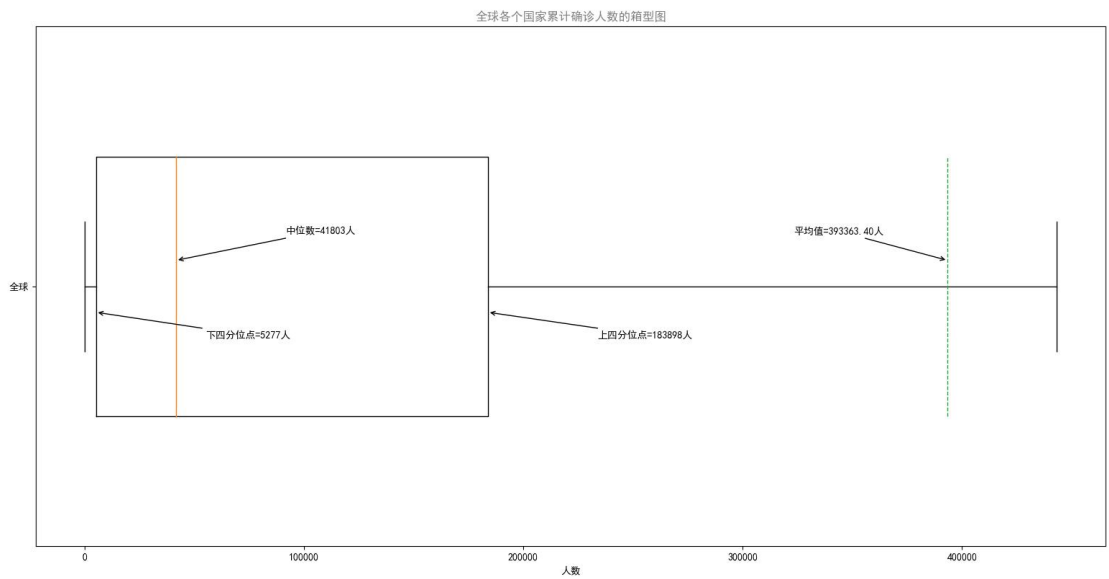


图 10 不含异常值的箱型图

8. 其他展示

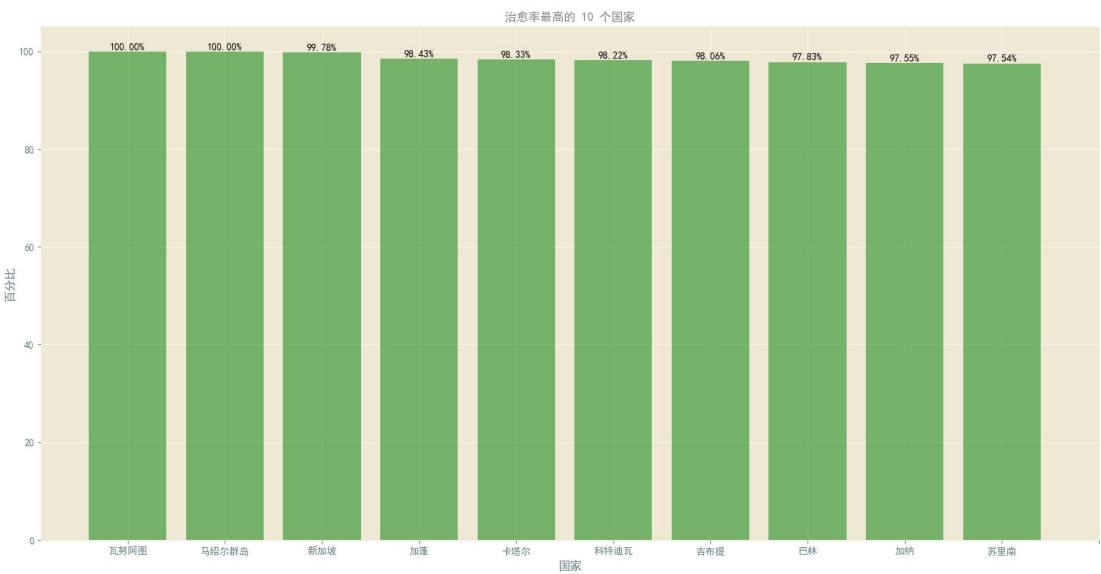


图 11 治愈率最高的 10 个国家

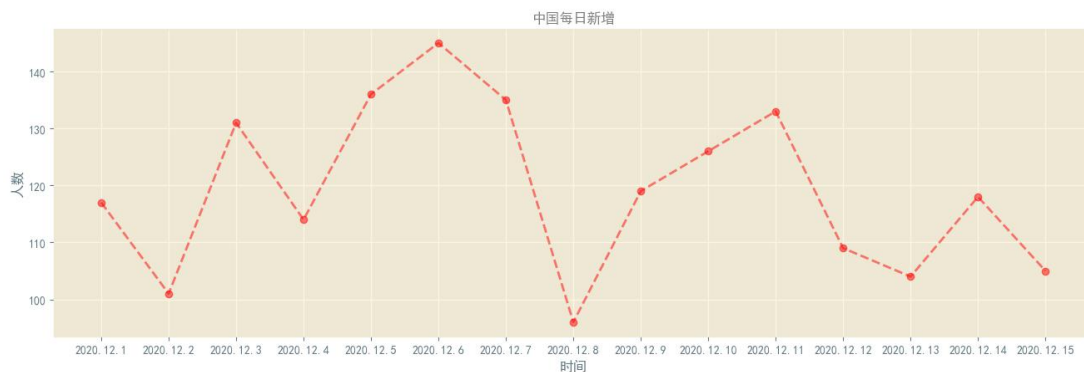


图 12 中国每日新增线图

#### 四. 根据以上数据，列出全世界应对新冠疫情最好的 10 个国家，并说明你的理由

对国家应对措施评分采用加权计分，满分 100，具体计算公式如下：

$$score = \text{治愈率} \times 15\% + (1 - \text{累计确诊比例}) \times 30\% + (1 - \text{死亡率}) \times 5\% + (1 - \text{新增确诊比例}) \times 50\%$$

对上式的解读如下：对疫情的应对应该体现在“对突发情况的处理以及进一步的对策”，体现在数据上就是累计确诊比例与新增确诊比例。故上式中这两项占比较大。

实际计算时，由于累计确诊比例，死亡率与新增确诊比例均较小，故还做了放大处理（即乘以合适的倍数）。

考虑到国与国之间人口数量差距较大，故采用百分比计算；同时由于部分国家人口过少，对数据影响较大，故不考虑总人口小于 5000 的国家。

最终得出分数最高的国家如下：

country	cure_per	total_per	death_per	new_per	people	score
泰国	93.005181	0.006138	1.413095	0.000024	6918	98.364474
尼日利亚	90.378063	0.037461	1.631368	0.000203	19587	97.813489
越南	88.873039	0.001453	2.496434	0.000004	9649	97.300627
菲律宾	92.702710	0.424222	1.950252	0.001265	10651	96.718083
巴基斯坦	87.159952	0.220729	2.009042	0.001501	20081	96.619189
中国	93.225356	0.006805	4.999633	0.000009	140005	96.597701
日本	81.776861	0.145658	1.387869	0.001830	12718	96.470036
埃塞俄比亚	81.220894	0.109911	1.542962	0.000482	10667	96.435279
印度	95.118908	0.731595	1.450703	0.002336	135405	96.297834
肯尼亚	79.791429	0.180677	1.730487	0.001143	5095	95.928904

图 13 应对较好的 10 个国家

观察上表，其中泰国优点在于累计确诊较少，治愈率也较高；越南这是累计确诊极少；中国优点在于治愈率较高，同时对新增确诊的控制较好。

总之，应对最好的 10 个国家是：泰国，尼日利亚，越南，菲律宾，巴基斯坦，中国，日本，埃塞俄比亚，印度，肯尼亚。



## 五. 针对全球累计确诊数，利用前 10 天采集到的数据做后 5 天的预测

考虑到数据量较小，预测效果可能不好，故除去上课介绍的线性回归预测方法外，还采用了两种对少量数据有较好预测效果的 GM(1,1)灰度预测与 k 阶差分预测。

三种预测在前 3 天都有较好的预测效果，误差基本低于 0.1%，其中 GM(1,1)灰度预测与 k 阶差分预测效果优于线性回归；而之后两天则出现误差的明显增加，误差达到约 1%：

GM(1,1)预测：				k阶差分预测：				一元线性回归预测：			
时间	预测	实际	误差	时间	预测	实际	误差	时间	预测	实际	误差
2020.12.11	70294657.27	70300495	-0.01%	2020.12.11	70288576.56	70300495	-0.02%	2020.12.11	70260855.80	70300495	-0.06%
2020.12.12	70954358.66	70985260	-0.04%	2020.12.12	70963576.82	70985260	-0.03%	2020.12.12	70892917.62	70985260	-0.13%
2020.12.13	71620251.23	71608902	0.02%	2020.12.13	71645298.58	71608902	0.05%	2020.12.13	71524979.44	71608902	-0.12%
2020.12.14	72292393.06	72984100	-0.95%	2020.12.14	72333515.37	72984100	-0.89%	2020.12.14	72157041.25	72984100	-1.13%
2020.12.15	72970842.80	73558656	-0.80%	2020.12.15	73028332.67	73558656	-0.72%	2020.12.15	72789103.07	73558656	-1.05%

图 14 三种预测比较

同时作图观察：

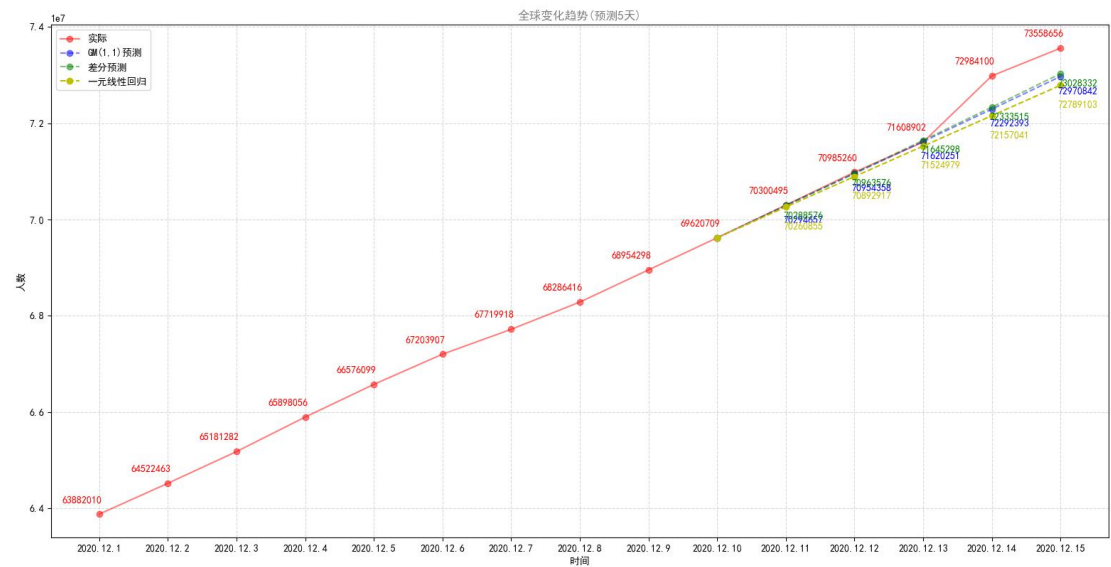


图 15 三种预测结果与实际的对比（全局趋势）

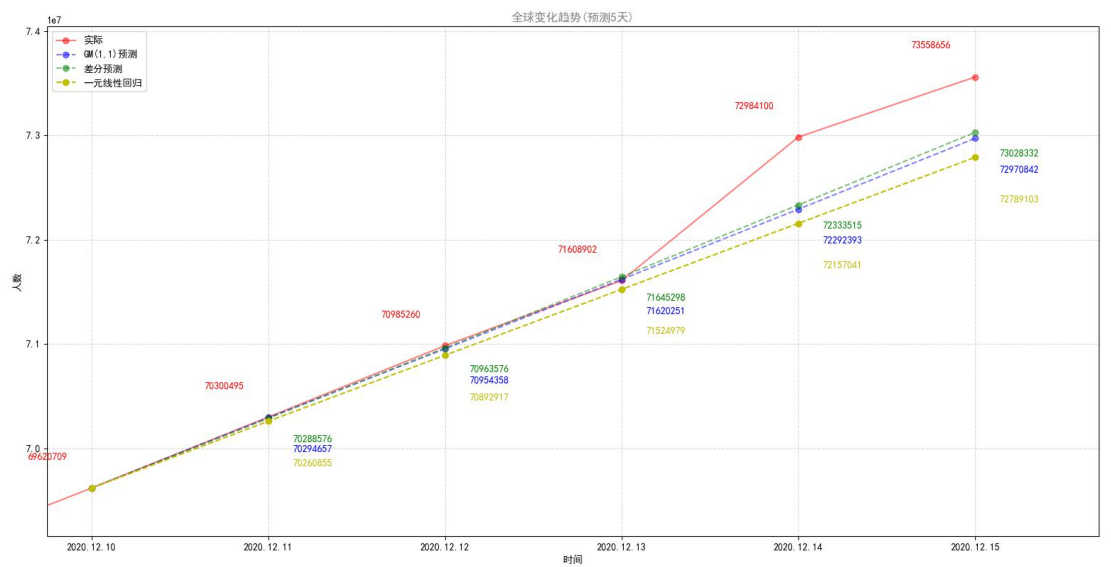


图 16 三种预测结果与实际的对比（放大预测部分）



从图中可以看到，12月14日累计确诊人数突然发生一次较大的上升，破坏了前13天的规律，导致预测不准。

考察12月14日的新增数据（图3，5），并未发现新增人数明显增加；查阅新闻也未发现当天有集中疫情爆发的国家。同时注意到，网站数据并不满足“后一天累计确诊=前一天累计确诊+后一条新增”，也不满足“各国每日新增之和等于全球每日新增数据”。于是有理由猜测，网站数据中“每日新增”与“累计确诊”来源于不同组织/机构发布；而且，与我国积极检测并上报的政策不同，其他国家的数据真实性与准确性无法保证。

考察腾讯的数据如下图，走向基本相同：

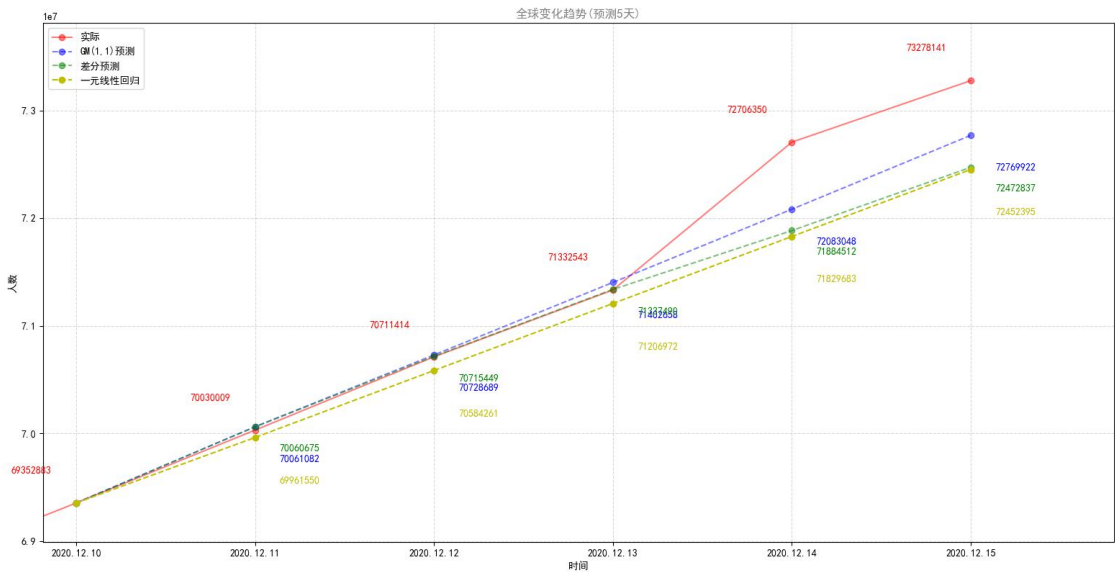


图 17 三种预测结果与实际的对比（腾讯数据）

最后，查看数据网站 <https://www.zq-ai.com/#/fe/xgfybigdata> 的曲线图（图18），观察到12月12日土耳其新增确诊异常提升，提升幅度与我绘制的线图基本一致，我爬取的同一网站数据在12日也有较高的每日新增数。至于为何累计确诊在14日才表现出异常，我认为数据源本身导致的不一致。

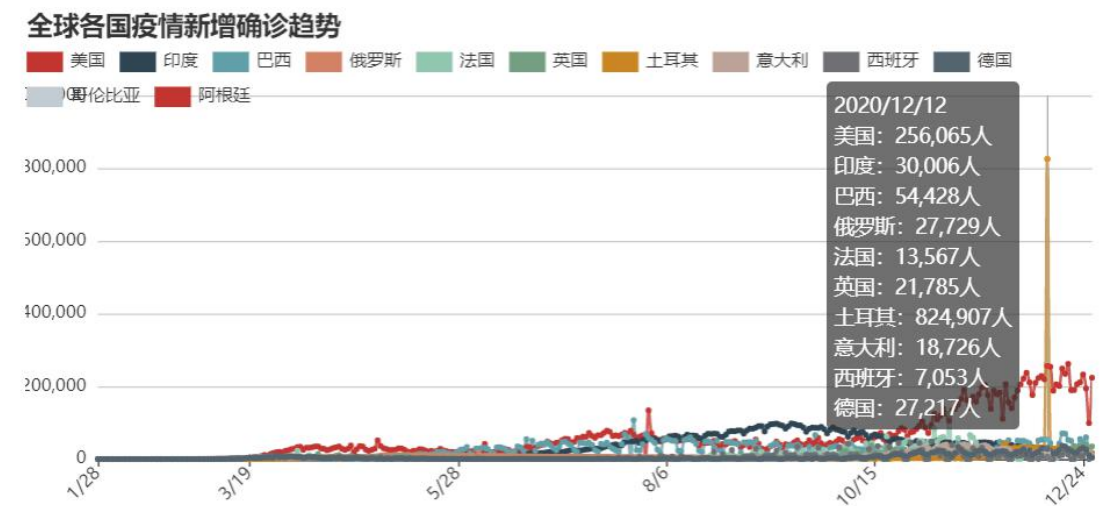


图 18 数据网站 <https://www.zq-ai.com/#/fe/xgfybigdata> 的曲线图

综合来看，三种数据预测方式都很好达到了预期效果，而12月14日的异常值属于现实世界中无法预测的突发事件。

## 六. 核心代码

### 1. 爬虫程序（使用 selenium）

spider:

```
1. import scrapy
2. from selenium import webdriver
3. from selenium.webdriver.chrome.options import Options
4. import datetime
5. from cov19.items import Cov19Item
6. from scrapy import Request
7.
8. class CovspiderSpider(scrapy.Spider):
9.     name = 'covSpider'
10.    start_urls = ['https://voice.baidu.com/act/newpneumonia/newpneumonia']
11.    all_urls = ['https://voice.baidu.com/act/newpneumonia/newpneumonia',
12.               'https://news.qq.com/zt2020/page/feiyang.htm#/global?nojump=1',
13.               'https://www.zq-ai.com/#/fe/xgfybigdata']
14.    showall = False
15.    site = 0 # 指示当前正在爬取哪个网站
16.
17.    def __init__(self):
18.        # 设定浏览器模式
19.        chrome_options = Options()
20.        chrome_options.add_argument('--headless') # 使用无头谷歌浏览器模式
21.        chrome_options.add_argument('--disable-gpu') # 禁用 GPU 加速
22.        chrome_options.add_argument('--no-sandbox')
23.        # 指定谷歌浏览器驱动路径
24.        self.driver = webdriver.Chrome(chrome_options=chrome_options,
25.                                       executable_path='C:/Users/87290/Desktop/PROGRAM/Python/chrome
26.                                       driver.exe')
27.        # 设定隐式等待时间 10s
28.        self.driver.implicitly_wait(10)
29.        # 输出提示信息
30.        print('Webdriver start init success!')
31.
32.    def parse(self, response):
33.        item = Cov19Item()
34.        print('*' * 40)
35.        print("开爬 开爬 fo~↑ fo~↑")
36.        if self.site == 0: # 如果是第一个网站(百度)
37.            print("# Baidu:")
38.            print("Get {} data".format(len(response.xpath('//*[id="foreignTable"]/table/tbody/tr/
39.            td/table/tbody/tr[*]'))))
```

```

38.         for each in response.xpath('//*[id="foreignTable"]/table/tbody/tr/td/table/tbody/tr[*]
        '):
39.             # 获取原始数据
40.             country = each.xpath('td[1]/a/div[1]/text()').extract()
41.             if len(country) == 0:
42.                 country = each.xpath('td[1]/div/text()').extract()
43.             country = country[0] # 获取国家名
44.             dailyNew = each.xpath('td[2]/text()').extract()[0] # 获取每日新增
45.             now = each.xpath('td[3]/text()').extract()[0] # 获取现有病例数
46.             total = each.xpath('td[4]/text()').extract()[0] # 获取累计确诊数
47.             cure = each.xpath('td[5]/text()').extract()[0] # 获取治愈数
48.             death = each.xpath('td[6]/text()').extract()[0] # 获取死亡数
49.
50.             # 装入 item
51.             item['year'] = datetime.datetime.now().year
52.             item['month'] = datetime.datetime.now().month
53.             item['day'] = datetime.datetime.now().day
54.             item['country'] = country.strip()
55.             item['dailyNew'] = dailyNew.strip()
56.             item['now'] = now.strip()
57.             item['total'] = total.strip()
58.             item['cure'] = cure.strip()
59.             item['death'] = death.strip()
60.             item['cure_rate'] = 'NA'
61.             item['death_rate'] = 'NA'
62.             item['site'] = 'Baidu'
63.             item['people'] = 'NA'
64.             yield item
65.
66.             # 提交下一个网址
67.             self.site += 1
68.             yield Request(self.all_urls[self.site])
69.
70.         elif self.site == 1: # 如果是第二个网站(腾讯)
71.             print("# Tencent:")
72.             print("Get {} data".format(len(response.xpath('//*[id="foreignWrapper"]/table/tbody[*]
        '))))
73.         for each in response.xpath('//*[id="foreignWrapper"]/table/tbody[*]'):
74.             # 获取原始数据
75.             country = each.xpath('tr/th/span/text()').extract()[0] # 获取国家名
76.             dailyNew = each.xpath('tr/td[1]/text()').extract()[0] # 获取每日新增
77.             total = each.xpath('tr/td[2]/text()').extract()[0] # 获取累计确诊数
78.             cure = each.xpath('tr/td[3]/text()').extract()[0] # 获取治愈数
79.             death = each.xpath('tr/td[4]/text()').extract()[0] # 获取死亡数

```

```

80.
81.         # 装入 item
82.         item['year'] = datetime.datetime.now().year
83.         item['month'] = datetime.datetime.now().month
84.         item['day'] = datetime.datetime.now().day
85.         item['country'] = country.strip()
86.         item['dailyNew'] = dailyNew.strip()
87.         item['now'] = 'NA'
88.         item['total'] = total.strip()
89.         item['cure'] = cure.strip()
90.         item['death'] = death.strip()
91.         item['cure_rate'] = 'NA'
92.         item['death_rate'] = 'NA'
93.         item['site'] = 'Tencent'
94.         item['people'] = 'NA'
95.         yield item
96.
97.         # 提交下一个网址
98.         self.site += 1
99.         yield Request(self.all_urls[self.site])
100.
101.     else:
102.         print("# AI Big Data:")
103.         print("Get {} data".format(len(response.xpath('//*[@id="app"]/div/section[2]/section/main/div/div[6]/div[3]/table/tbody/tr[*]'))))
104.         for each in response.xpath('//*[@id="app"]/div/section[2]/section/main/div/div[6]/div[3]/table/tbody/tr[*]'):
105.             # 获取原始数据
106.             country = each.xpath('td[1]/div/a/span/text()').extract()[0]
107.             dailyNew = each.xpath('td[2]/div/span/text()').extract()[0]
108.             now = each.xpath('td[4]/div/span/text()').extract()[0]
109.             total = each.xpath('td[3]/div/span/text()').extract()[0]
110.             cure = each.xpath('td[6]/div/span/text()').extract()[0]
111.             death = each.xpath('td[9]/div/span/text()').extract()[0]
112.             cure_rate = each.xpath('td[7]/div/span/text()').extract()[0]
113.             death_rate = each.xpath('td[10]/div/span/text()').extract()[0]
114.             people = each.xpath('td[12]/div/span/text()').extract()[0].replace(',', '')
115.
116.             # 装入 item
117.             item['year'] = datetime.datetime.now().year
118.             item['month'] = datetime.datetime.now().month
119.             item['day'] = datetime.datetime.now().day
120.             item['country'] = country.strip()
121.             item['dailyNew'] = dailyNew.strip()

```

```

122.         item['now'] = now.strip()
123.         item['total'] = total.strip()
124.         item['cure'] = cure.strip()
125.         item['death'] = death.strip()
126.         item['cure_rate'] = cure_rate.strip()
127.         item['death_rate'] = death_rate.strip()
128.         item['site'] = 'AI'
129.         item['people'] = people.strip()
130.         yield item
131.         pass
132.
133.     print('*' * 40)
134.     pass

```

middlewares:

```

1.     import scrapy
2.     import time
3.
4.     class Cov19Middleware(object):
5.         def __init__(self):
6.             pass
7.         def process_request(self, request, spider):
8.             wd = spider.driver
9.             # 浏览器加载链接
10.            wd.get(request.url)
11.
12.            # 数据完整性保证
13.            if spider.site == 0: # 如果是百度, 要展开全部数据
14.                if spider.showall == False: # 如果还没展开全部数据, 就通过 click 点击按钮使其展开
15.                    print("Try to show all details...")
16.                    all = wd.find_element_by_xpath('//*[@id="foreignTable"]/div/span')
17.                    all.click()
18.                    spider.showall = True
19.                else: # 否则等待两秒确保加载完毕
20.                    print('continue')
21.                    time.sleep(5)
22.
23.            # 打包返回 response
24.            html = spider.driver.page_source
25.            return scrapy.http.HtmlResponse(url=spider.driver.current_url, body=html.encode('utf-8'),
                encoding='utf-8',
26.
                request=request)

```

pipelines:

```

1.     import csv
2.

```

```

3. class Cov19Pipeline(object):
4.     all_path = []
5.     def open_spider(self, spider):
6.         self.all_path.append("cov19_Baidu.csv")
7.         self.all_path.append("cov19_Tencent.csv")
8.         self.all_path.append("cov19_AI.csv")
9.         # for path in self.all_path:
10.        #     # 写表头
11.        #     with open(path, 'w', encoding="utf-8", newline='') as f:
12.        #         csv_write = csv.writer(f)
13.        #         csv_head = ['year', 'month', 'day', 'country', 'dailyNew', 'now', 'total', 'cure',
14.        #                     'cure_rate', 'death_rate', 'site', 'people']
15.        #         csv_write.writerow(csv_head)
16.
17.     def process_item(self, item, spider):
18.         # 确定往哪张表里添加数据
19.         path = self.all_path[spider.site]
20.         # 写数据
21.         with open(path, 'a', encoding="utf-8", newline='') as f:
22.             csv_write = csv.writer(f)
23.             csv_write.writerow(item.values())
24.
25.         return item
26.
27.     def close_spider(self, spider):
28.         spider.driver.quit()
29.         print('Chrome closed!')
30.         print('Spider closed!')

```

## 2. 各国应对措施打分代码

```

1. import pandas as pd
2.
3. # 读取数据(直接前面的数据)
4. cure = pd.read_csv("cure_1.csv")
5. cure.columns = pd.Series(['country', 'cure_per']) # 修改列名
6. total = pd.read_csv("d_1.csv")
7. total.columns = pd.Series(['country', 'total_per']) # 修改列名
8. death = pd.read_csv("e_1.csv")
9. death.columns = pd.Series(['country', 'death_per']) # 修改列名
10. new = pd.read_csv("c1_1.csv")
11. new.columns = pd.Series(['country', 'new_per']) # 修改列名
12. people = pd.read_csv("people.csv")
13. new.columns = pd.Series(['country', 'people']) # 修改列名
14.

```

```

15. key = pd.concat([cure, total, death, new, people], axis=1, join='outer', ignore_index=True) # 合并表
16. key.drop([2, 4, 6, 8], inplace=True, axis=1) # 删除冗余列
17. key.columns = pd.Series(['country', 'cure_per', 'total_per', 'death_per', 'new_per', 'people']) # 修改列名
18. key.drop(index=(key.loc[(key['country']=='钻石号邮轮')].index), inplace=True) # 删除特殊数据
19. key.drop(index=(key.loc[(key['country']=='日本本土')].index), inplace=True) # 删除特殊数据
20. key.drop(index=(key.loc[(key['death_per']==0)].index), inplace=True) # 删除死亡人数为0的行
21. key.drop(index=(key.loc[(key['people']<5000)].index), inplace=True) # 删除人口小于5000的国家
22.
23. key['score'] = (key['cure_per'] * 13 + (100 - key['total_per'] * 10) * 32 + (100 - key['death_per'] * 10) * 5 + (100 - key['new_per'] * 1) * 50) / 100
24. key.sort_values(by='score', ascending=False, inplace=True) # 排序
25. print(key)
26. key.to_csv('score.csv')

```

### 3. 数据展示代码

#### a. 15 天中，全球新冠疫情的总体变化趋势

```

1. import pandas as pd
2. import numpy as np
3. import matplotlib
4. import matplotlib.pyplot as plt
5.
6. # 15 天中，全球新冠疫情的总体变化趋势
7. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
8. plt.style.use('Solarize_Light2')
9.
10. fig = plt.figure()
11. # fig, ax = plt.subplots()
12. ax1 = fig.add_subplot(211)
13. ax2 = fig.add_subplot(212)
14. csv_path = "cov19_Tencent.csv"
15. data = pd.read_csv(csv_path)
16. data[(data['dailyNew'] < 0)]['dailyNew'] = np.nan # 将异常数据置为0
17.
18. key = [data.groupby(['year', 'month', 'day'])['total'].sum()] # 求得每日全球累计确诊数之和
19. key.append(data.groupby(['year', 'month', 'day'])['dailyNew'].sum())
20.
21. keyy = pd.concat(key, axis=1)
22.
23. print(keyy)
24.
25. keyy.to_csv('a_1.csv') # 储存每日全球累计确诊数之和
26. keyy = pd.read_csv('a_1.csv')
27. print(keyy)

```



```

28.
29. x1 = np.arange(1, keyy.shape[0] + 1, 1) # 图形 1x 数据
30. y1 = (keyy['total']).tolist() # 图形 1y 数据
31. ax1.plot(x1, y1, marker='o', linestyle='--', alpha=0.5, color='g') # 作图
32.
33. x2 = np.arange(1, keyy.shape[0] + 1, 1) # 图形 2x 数据
34. y2 = (keyy['dailyNew']).tolist() # 图形 2y 数据
35. ax2.plot(x2, y2, marker='o', linestyle='--', alpha=0.5, color='r') # 作图
36.
37. xticks = [] # 构造 x 标签
38. for i in range(keyy.shape[0]):
39.     tmplist = keyy.iloc[i].tolist()
40.     xticks.append(str(tmplist[0])+'.'+str(tmplist[1])+'.'+str(tmplist[2]))
41.     pass
42. ax1.set_xticks(range(1, keyy.shape[0] + 1)) # 设置 x 轴刻度
43. ax1.set_xticklabels(xticks) # 设置 x 轴刻度标签
44. ax1.set_xlabel("时间") # 设置 x 轴上的标签
45. ax1.set_ylabel("人数") # 设置 y 轴上的标签
46. ax1.set_title("全球变化趋势(累计)", fontsize=12, color='grey') # 设置标题
47.
48. ax2.set_xticks(range(1, keyy.shape[0] + 1)) # 设置 x 轴刻度
49. ax2.set_xticklabels(xticks) # 设置 x 轴刻度标签
50. ax2.set_xlabel("时间") # 设置 x 轴上的标签
51. ax2.set_ylabel("人数") # 设置 y 轴上的标签
52. ax2.set_title("全球变化趋势(新增)", fontsize=12, color='grey') # 设置标题
53.
54. for a, b in zip(x1, y1):
55.     ax1.text(a - 0.25, b, '%d' % b, ha='center', va='center', fontsize=10) # 数值点旁添加具体值
56. for a, b in zip(x2, y2):
57.     ax2.text(a - 0.2, b, '%d' % b, ha='center', va='center', fontsize=10) # 数值点旁添加具体值
58.
59. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
60. plt.show() # 显示

```

b. 累计确诊数排名前 20 的国家名称及其数量

```

1. import pandas as pd
2. import numpy as np
3. import matplotlib
4. import matplotlib.pyplot as plt
5.
6. # 累计确诊数排名前 20 的国家名称及其数量
7. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
8. plt.style.use('Solarize_Light2')
9.
10. fig, ax = plt.subplots()

```

```

11. csv_path = "cov19_AI.csv"
12. data = pd.read_csv(csv_path)
13.
14. front = 20 # 取前 20 个国家
15.
16. key = [data.groupby(['country'])['total'].max()] # 每个国家累计确诊数
17. keyy = pd.concat(key, axis=1)
18. keyy.sort_values(by='total', ascending=False, inplace=True)
19. keyy.to_csv('b_1.csv') # 储存每日全球累计确诊数之和
20. keyy = pd.read_csv('b_1.csv')
21. print(keyy)
22.
23. x = np.arange(1, front + 1, 1) # 图形 x 数据
24. y = [each for each in (keyy['total'] - keyy.iloc[front - 1][1] + 100000).tolist()[0:front]] # 图形 y 数据
25. ax.bar(x, y, alpha=0.5, color='g') # 作条形图
26.
27. xticks = [] # 构造 x 标签
28. for i in range(front):
29.     tmplist = keyy.iloc[i].tolist()
30.     xticks.append(tmplist[0])
31.     pass
32. ax.set_xticks(range(1, front + 2)) # 设置 x 轴刻度
33. ax.set_xticklabels(xticks) # 设置 x 轴刻度标签
34. ax.set_xlabel("国家") # 设置 x 轴上的标签
35. ax.set_ylabel("人数(为观察方便, 已将全部数据减去{})".format(keyy.iloc[front - 1][1] + 100000)) # 设置 y 轴上的标签
36. ax.set_title("累计确诊前{}国".format(front), fontsize=12, color='grey') # 设置标题
37.
38. for a, b in zip(x, y):
39.     ax.text(a, b + keyy.iloc[0][1] // 70, '%d' % (b + keyy.iloc[front - 1][1] - 1),
40.             ha='center', va='center', fontsize=10) # 数值点旁添加具体值
41.
42. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
43. plt.show() # 显示

```

c. 15 天中, 每日新增确诊数累计排名前 10 个国家的每日新增确诊数据的曲线图

```

1. import pandas as pd
2. import numpy as np
3. import matplotlib
4. import matplotlib.pyplot as plt
5.
6. # 15 天中, 每日新增确诊数累计排名前 10 个国家的每日新增确诊数据的曲线图
7. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
8. plt.style.use('Solarize_Light2')

```

```

9.
10. fig, ax = plt.subplots()
11. csv_path = "cov19_AI.csv"
12. data = pd.read_csv(csv_path)
13.
14. front = 10 # 取前 10 个国家
15.
16. key = [data.groupby(['country'])['dailyNew'].sum()] # 求得每个国家累计确诊数之和
17. keyy = pd.concat(key, axis=1)
18. keyy.sort_values(by='dailyNew', ascending=False, inplace=True)
19. keyy.to_csv('c_1.csv') # 储存
20. keyy = pd.read_csv('c_1.csv')
21. front_ten = [] # 提取前 10 个国家名字
22. for i in range(front):
23.     tmp_list = keyy.iloc[i].tolist()
24.     front_ten.append(tmp_list[0])
25.     pass
26. print("15 天中，每日新增确诊数累计排名前 10 的国家")
27. print(front_ten)
28. print()
29.
30. # 画 10 条曲线
31. for i in range(len(front_ten)):
32.     tmp_data = data[(data['country'] == front_ten[i])]
33.     x = np.arange(1, tmp_data.shape[0] + 1, 1) # 图形 x 数据
34.     y = (tmp_data['dailyNew']).tolist() # 图形 y 数据
35.     plt.plot(x, y)
36.     if i == 0:
37.         for a, b in zip(x, y):
38.             ax.text(a, b - 7000, '%d' % (b + keyy.iloc[front - 1][1] - 1),
39.                     ha='center', va='center', fontsize=10) # 数值点旁添加具体值
40. xticks = [] # 构造 x 标签
41. for i in range(tmp_data.shape[0]):
42.     tmp_list = tmp_data.iloc[i].tolist()
43.     xticks.append(str(tmp_list[0]) + '.' + str(tmp_list[1]) + '.' + str(tmp_list[2]))
44.     pass
45. ax.set_xticks(range(1, tmp_data.shape[0] + 1)) # 设置 x 轴刻度
46. print(tmp_data)
47.
48. ax.legend(front_ten, loc='upper left') # 设置图例
49. ax.set_xticklabels(xticks) # 设置 x 轴刻度标签
50. ax.set_xlabel("时间") # 设置 x 轴上的标签
51. ax.set_ylabel("人数") # 设置 y 轴上的标签
52. ax.set_title("每日新增确诊数据", fontsize=12, color='grey') # 设置标题

```

53.

```
54. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
```

```
55. plt.show() # 显示
```

d. 累计确诊人数占国家总人口比例最高的 10 个国家

```
1. import pandas as pd
```

```
2. import numpy as np
```

```
3. import matplotlib
```

```
4. import matplotlib.pyplot as plt
```

```
5.
```

```
6. # 累计确诊人数占国家总人口比例最高的 10 个国家
```

```
7. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
```

```
8. plt.style.use('Solarize_Light2')
```

```
9.
```

```
10. fig, ax = plt.subplots()
```

```
11. csv_path = "cov19_AI.csv"
```

```
12.
```

```
13. front = 10 # 取前 10 个国家
```

```
14.
```

```
15. data = pd.read_csv(csv_path)
```

```
16. data['people'] = data['people'].replace('-', 'NaN') # 数据清洗: 用 NaN 代替-
```

```
17. data['people'] = data['people'].astype('float') # 数据转化为 float
```

```
18.
```

```
19. key = [data.groupby(['country'])['total'].max() / (data.groupby(['country'])['people'].max() * 100)  
] # 每个国家百分比
```

```
20. keyy = pd.concat(key, axis=1) # 转化为 df
```

```
21. keyy.to_csv('d_1.csv') # 储存
```

```
22. keyy = pd.read_csv('d_1.csv')
```

```
23. keyy.columns = pd.Series(['country', 'per']) # 修改列名
```

```
24. keyy.drop(index=(keyy.loc[(keyy['country']=='钻石号邮轮')].index), inplace=True) # 删除特殊数据
```

```
25. keyy.sort_values(by='per', ascending=False, inplace=True) # 排序
```

```
26. print(keyy)
```

```
27.
```

```
28. x = np.arange(1, front + 1, 1) # 图形 x 数据
```

```
29. y = [each for each in keyy['per'].tolist()[0:front]] # 图形 y 数据
```

```
30. ax.bar(x, y, alpha=0.5, color='g') # 作条形图
```

```
31.
```

```
32. xticks = [] # 构造 x 标签
```

```
33. for i in range(front):
```

```
34.     tmplist = keyy.iloc[i].tolist()
```

```
35.     xticks.append(tmplist[0])
```

```
36.     pass
```

```
37. ax.set_xticks(range(1, front + 2)) # 设置 x 轴刻度
```

```
38. ax.set_xticklabels(xticks) # 设置 x 轴刻度标签
```

```
39. ax.set_xlabel("国家") # 设置 x 轴上的标签
```

```

40. ax.set_ylabel("百分比") # 设置 y 轴上的标签
41. ax.set_title("累计确诊人数占国家总人口比例", fontsize=12, color='grey') # 设置标题
42. print(xticks)
43.
44. for a, b in zip(x, y):
45.     ax.text(a, b + 0.1, '{:.2f}%'.format(b),
46.             ha='center', va='center', fontsize=10) # 数值点旁添加具体值
47.
48. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
49. plt.show() # 显示

```

e. 死亡率（累计死亡人数/累计确诊人数）最低的 10 个国家

```

1. import pandas as pd
2. import numpy as np
3. import matplotlib
4. import matplotlib.pyplot as plt
5. # 死亡率
6.
7. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
8. plt.style.use('Solarize_Light2')
9.
10. fig, ax = plt.subplots()
11. csv_path = "cov19_AI.csv"
12.
13. front = 10 # 取前 10 个国家
14.
15. data = pd.read_csv(csv_path)
16. data['death'] = data['death'].astype('float') # 数据转化为 float
17. data['total'] = data['total'].astype('float') # 数据转化为 float
18.
19. key = [data.groupby(['country'])['death'].max() / data.groupby(['country'])['total'].max() * 100]
    # 每个国家百分比
20. keyy = pd.concat(key, axis=1) # 转化为 df
21. keyy.to_csv('e_1.csv') # 储存
22. keyy = pd.read_csv('e_1.csv')
23. keyy.columns = pd.Series(['country', 'death_per']) # 修改列名
24. keyy.drop(index=(keyy.loc[(keyy['country']=='钻石号邮轮').index], inplace=True) # 删除特殊数据
25. keyy.drop(index=(keyy.loc[(keyy['death_per']==0)].index), inplace=True) # 删除死亡人数为 0 的行
26. keyy.sort_values(by='death_per', ascending=True, inplace=True) # 排序
27. print(keyy)
28.
29. x = np.arange(1, front + 1, 1) # 图形 x 数据
30. y = [each for each in keyy['death_per'].tolist()[0:front]] # 图形 y 数据
31. ax.bar(x, y, alpha=0.5, color='g') # 作条形图
32.

```

```

33. xticks = [] # 构造 x 标签
34. for i in range(front):
35.     tmplist = keyy.iloc[i].tolist()
36.     xticks.append(tmplist[0])
37.     pass
38. ax.set_xticks(range(1, front + 2)) # 设置 x 轴刻度
39. ax.set_xticklabels(xticks) # 设置 x 轴刻度标签
40. ax.set_xlabel("国家") # 设置 x 轴上的标签
41. ax.set_ylabel("百分比") # 设置 y 轴上的标签
42. ax.set_title("死亡率最低的 10 个国家", fontsize=12, color='grey') # 设置标题
43. print(xticks)
44.
45. for a, b in zip(x, y):
46.     ax.text(a, b + 0.01, '{:.2f}%'.format(b),
47.             ha='center', va='center', fontsize=10) # 数值点旁添加具体值
48.
49. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
50. plt.show() # 显示

```

f. 用饼图展示各个国家的累计确诊人数的比例

```

1. import pandas as pd
2. import numpy as np
3. import matplotlib
4. import matplotlib.pyplot as plt
5. from matplotlib import font_manager as fm
6.
7. # 用饼图展示各个国家的累计确诊人数的比例（你爬取的所有国家，数据较小的国家可以合并处理）
8. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
9.
10. fig, ax = plt.subplots()
11. csv_path = "b_1.csv"
12. keyy = pd.read_csv(csv_path) # 直接利用第 b 问的数据
13.
14. border = 600000 # 合并小于 border 国家 为 "其他"
15. borderBySort = keyy.iloc[9][1] # 或者合并 10 名以后的国家 为 "其他"
16.
17. keyy.loc[keyy.shape[0]] = ["其他", 0] # 添加 "其他" 列，初始化为 0
18.
19. drop_lines = [] # 记录要删除的行
20.
21. for i in range(keyy['country'].count() - 1):
22.     if keyy.iloc[i][1] < borderBySort:
23.         keyy.iloc[keyy.loc[(keyy['country'] == '其他').index, 1] += keyy.iloc[i][1] # 修改 "其他"
            " 行
24.         drop_lines.append(i) # 记录要删除的行

```

```

25.         pass
26.     pass
27.     keyy.drop(index=drop_lines, inplace=True) # 删除数据
28.     keyy = keyy.reset_index(drop=True) # 重设索引
29.
30.     print(keyy)
31.
32.     # 作图
33.     size = [] # 各国家的人数
34.     label_list = [] # 各部分标签
35.     for i in range(keyy.shape[0]):
36.         tmplist = keyy.iloc[i].tolist()
37.         size.append(tmplist[1])
38.         label_list.append(tmplist[0])
39.     pass
40.
41.     explode = np.zeros(keyy.shape[0]) # 各部分的突出显示比例(法 1)
42.     explode[0] = 0.1
43.     # explode = np.linspace(0.2, 0.05, keyy.shape[0]) # 各部分的突出显示比例(法 2)
44.     patches, texts, autotexts = plt.pie(size, explode=explode, labels=label_list,
45.                                         labeldistance=1.1, autopct="%1.1f%%", shadow=False, startangle=
46.                                         0, pctdistance=0.9)
47.     # 调整字体
48.     proptease = fm.FontProperties()
49.     proptease.set_size('x-large')
50.     plt.setp(texts, fontproperties=proptease)
51.     plt.setp(autotexts, fontproperties=proptease)
52.     ax.set_title('各个国家的累计确诊人数的比例(仅展示前 10 个国家, 剩下的归于"其他")', fontsize=16) # 设置标
53.     # 窗口最大化
54.     plt.get_current_fig_manager().window.showMaximized()
55.     plt.show()

```

g. 展示全球各个国家累计确诊人数的箱型图，要有平均值

```

1.     import pandas as pd
2.     import matplotlib
3.     import matplotlib.pyplot as plt
4.     from util import util
5.
6.     # 展示全球各个国家累计确诊人数的箱型图，要有平均值
7.     matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
8.
9.     fig, ax = plt.subplots()
10.    csv_path = "b_1.csv"
11.    keyy = pd.read_csv(csv_path) # 直接利用第 b 问的数据
12.    avee = keyy['total'].mean()

```



```

13. maxx = keyy['total'].max()
14. medd = keyy['total'].median()
15. upperQuartile = keyy['total'].quantile(0.75)
16. lowerQuartile = keyy['total'].quantile(0.25)
17.
18. x = keyy['total'].tolist() # 图形数据
19. f = plt.boxplot(x, showmeans=True, vert=False, meanline=True, widths=0.5, showfliers=False)
20.
21. ax.set_xlabel("人数") # 设置x轴上的标签
22. ax.set_yticklabels(['全球']) # 设置y轴刻度标签
23. ax.set_title("全球各个国家累计确诊人数的箱型图", fontsize=12, color='grey') # 设置标题
24.
25. # 设置关键点注释
26. ax.annotate("平均值={:.2f}人".format(avee), xy=(avee, 1.05),
27.             xytext=(avee - 70000, 1.1), arrowprops=dict(arrowstyle="->"))
28. ax.annotate("最大值={:.0f}人".format(maxx), xy=(maxx, 1),
29.             xytext=(maxx - 50000, 1.1), arrowprops=dict(arrowstyle="->"))
30. ax.annotate("中位数={:.0f}人".format(medd), xy=(medd, 1.05),
31.             xytext=(medd + 50000, 1.1), arrowprops=dict(arrowstyle="->"))
32. ax.annotate("上四分位点={:.0f}人".format(upperQuartile), xy=(upperQuartile, 0.95),
33.             xytext=(upperQuartile + 50000, 0.9), arrowprops=dict(arrowstyle="->"))
34. ax.annotate("下四分位点={:.0f}人".format(lowerQuartile), xy=(lowerQuartile, 0.95),
35.             xytext=(lowerQuartile + 50000, 0.9), arrowprops=dict(arrowstyle="->"))
36.
37. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
38. plt.show()

```

#### h. 其他展示

```

1. import pandas as pd
2. import numpy as np
3. import matplotlib
4. import matplotlib.pyplot as plt
5.
6. # 中国
7. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
8. plt.style.use('Solarize_Light2')
9.
10. fig, ax = plt.subplots()
11. csv_path = "cov19_AI.csv"
12. data = pd.read_csv(csv_path)
13.
14. keyyy = data[(data['country'] == '中国')][['year', 'month', 'day', 'dailyNew']]
15.
16. x = np.arange(1, keyyy.shape[0] + 1, 1) # 图形 2x 数据
17. y = (keyyy['dailyNew']).tolist() # 图形 2y 数据

```

```

18. ax.plot(x, y, marker='o', linestyle='--', alpha=0.5, color='r') # 作图
19.
20. xticks = [] # 构造 x 标签
21. for i in range(keyyy.shape[0]):
22.     tmplist = keyyy.iloc[i].tolist()
23.     xticks.append(str(tmplist[0])+'.'+str(tmplist[1])+'.'+str(tmplist[2]))
24.     pass
25.
26. ax.set_xticks(range(1, keyyy.shape[0] + 1)) # 设置 x 轴刻度
27. ax.set_xticklabels(xticks) # 设置 x 轴刻度标签
28. ax.set_xlabel("时间") # 设置 x 轴上的标签
29. ax.set_ylabel("人数") # 设置 y 轴上的标签
30. ax.set_title("中国每日新增", fontsize=12, color='grey') # 设置标题
31. plt.show()

```

```

1. import pandas as pd
2. import numpy as np
3. import matplotlib
4. import matplotlib.pyplot as plt
5.
6. # 治愈
7. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
8. plt.style.use('Solarize_Light2')
9.
10. fig, ax = plt.subplots()
11. csv_path = "cov19_AI.csv"
12.
13. front = 10 # 取前 10 个国家
14.
15. data = pd.read_csv(csv_path)
16. data['death'] = data['cure'].astype('float') # 数据转化为 float
17. data['total'] = data['total'].astype('float') # 数据转化为 float
18.
19. key = [data.groupby(['country'])['cure'].max() / data.groupby(['country'])['total'].max() * 100]
    # 每个国家百分比
20. keyy = pd.concat(key, axis=1) # 转化为 df
21. keyy.to_csv('cure_1.csv') # 储存
22. keyy = pd.read_csv('cure_1.csv')
23. keyy.columns = pd.Series(['country', 'cure_per']) # 修改列名
24. keyy.drop(index=(keyy.loc[(keyy['country']=='钻石号邮轮').index], inplace=True) # 删除特殊数据
25. keyy.drop(index=(keyy.loc[(keyy['cure_per']==0).index], inplace=True) # 删除死亡人数为 0 的行
26. keyy.sort_values(by='cure_per', ascending=False, inplace=True) # 排序
27. print(keyy)
28.

```

```

29. x = np.arange(1, front + 1, 1) # 图形 x 数据
30. y = [each for each in keyy['cure_per'].tolist()[0:front]] # 图形 y 数据
31. ax.bar(x, y, alpha=0.5, color='g') # 作条形图
32.
33. xticks = [] # 构造 x 标签
34. for i in range(front):
35.     tmplist = keyy.iloc[i].tolist()
36.     xticks.append(tmplist[0])
37.     pass
38. ax.set_xticks(range(1, front + 2)) # 设置 x 轴刻度
39. ax.set_xticklabels(xticks) # 设置 x 轴刻度标签
40. ax.set_xlabel("国家") # 设置 x 轴上的标签
41. ax.set_ylabel("百分比") # 设置 y 轴上的标签
42. ax.set_title("治愈率最高的 10 个国家", fontsize=12, color='grey') # 设置标题
43.
44. for a, b in zip(x, y):
45.     ax.text(a, b + 0.9, '{:.2f}%'.format(b),
46.             ha='center', va='center', fontsize=10) # 数值点旁添加具体值
47.
48. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
49. plt.show() # 显示

```

#### 4. 数据预测代码

```

1. import pandas as pd
2. import matplotlib
3. import matplotlib.pyplot as plt
4. import numpy as np
5. from scipy.optimize import nnls
6. import math
7. from sklearn.linear_model import LinearRegression
8.
9.
10. # 图形界面预设以及生成画布与坐标轴
11. matplotlib.rcParams['font.sans-serif'] = ['SimHei'] # 设置对中文的支持
12. fig, ax = plt.subplots()
13.
14. # 读取数据
15. csv_path = "a_1.csv"
16. keyy = pd.read_csv(csv_path) # 直接利用第 a 问的数据
17.
18. # 进一步设定初始数据
19. preDays = 5 # 预测的天数
20. real = keyy['total'].values.tolist() # 取得真实数据
21. history = keyy['total'].values.tolist() # 取得历史数据

```

```

22. for i in range(preDays):
23.     history.pop() # 删除后几天的数据
24.
25. # ----- #
26. # GM(1,1)预测
27.
28. pre1 = [] # 存放预测的数据
29. n = len(history) # 用于训练的数据量
30. pre1.append(history[-1]) # 便于画图，添加一项真实数据
31.
32. X0 = np.array(history)
33. X1 = np.array([sum(history[0:i+1]) for i in range(n)]) # 对数据作一次累加
34.
35. # 构造数据矩阵 B 及数据向量 Y
36. B = np.zeros([n-1, 2])
37. Y = np.zeros([n-1, 1])
38. for i in range(0, n-1):
39.     B[i, 0] = -0.5*(X1[i] + X1[i+1])
40.     B[i, 1] = 1
41.     Y[i, 0] = X0[i+1]
42.
43. # 计算 GM(1,1)微分方程的参数 a 和 b
44. U = np.linalg.inv(B.T.dot(B)).dot(B.T).dot(Y)
45. a = U[0][0]
46. b = U[1][0]
47.
48. # 预测
49. f = np.zeros(preDays)
50. print('GM(1,1)预测:')
51. print('时间      预测      实际      误差')
52. for i in range(0, preDays):
53.     tmplist = keyy.iloc[keyy.shape[0] - preDays + i].tolist()
54.     date = str(tmplist[0]) + '.' + str(tmplist[1]) + '.' + str(tmplist[2]) + '\t'
55.     f[i] = (X0[0] - b / a) * (1 - math.exp(a)) * math.exp(-a * (i + n)) # 此次预测的值
56.     pre1.append(int(f[i]))
57.     print(date, '{:.2f}'.format(f[i]), '\t', real[keyy.shape[0] - preDays + i], '\t',
58.           '{:.2f}%'.format((f[i] - real[keyy.shape[0] - preDays + i]) / real[keyy.shape[0] - preDa
59. ys + i] * 100))
60. # ----- #
61. # k阶差分预测
62. pre2 = []
63. pre2.append(history[-1])
64.

```

```

65. k = 8 # k 阶差分预测
66.
67. n = len(history)
68. x = [history[k - i: n-i] for i in range(1, k+1)]
69. x.append([1 for j in range(n-k)])
70. x = np.array(x).T
71. y = np.array(history[k: n]).T
72.
73. paras = nnls(x, y)[0]
74.
75. sq = history.copy()
76. j = n
77. print('k 阶差分预测:')
78. print('时间          预测          实际          误差')
79. for i in range(preDays):
80.     tmplist = keyy.iloc[keyy.shape[0] - preDays + i].tolist()
81.     date = str(tmplist[0]) + '.' + str(tmplist[1]) + '.' + str(tmplist[2]) + '\t'
82.     temp = 0
83.     for j in range(k):
84.         temp += paras[j] * sq[len(sq)-j-1]
85.     temp += paras[k]
86.     pre2.append(int(temp))
87.     sq.append(temp)
88.     # print(temp)
89.     print(date, '{:.2f}'.format(temp), '\t', real[keyy.shape[0] - preDays + i], '\t',
90.           '{:.2f}%'.format((temp - real[keyy.shape[0] - preDays + i]) / real[keyy.shape[0] - preDays + i] * 100))
91.
92.
93. # ----- #
94. # 一元线性回归
95. pre3 = []
96. pre3.append(history[-1])
97.
98. Model = LinearRegression()
99. Model.fit(np.arange(1, keyy.shape[0] + 1 - preDays, 1).reshape(-1, 1), history)
100.
101. print('一元线性回归预测:')
102. print('时间          预测          实际          误差')
103. for i in range(0, preDays):
104.     tmplist = keyy.iloc[keyy.shape[0] - preDays + i].tolist()
105.     date = str(tmplist[0]) + '.' + str(tmplist[1]) + '.' + str(tmplist[2]) + '\t'
106.     tmppre = Model.predict([[keyy.shape[0] - preDays + i + 1]]).tolist()[0]
107.     pre3.append(int(tmppre))

```

```

108.     print(date, '{:.2f}'.format(tmppre), '\t', real[keyy.shape[0] - preDays + i], '\t',
109.           '{:.2f}%'.format((tmppre - real[keyy.shape[0] - preDays + i]) / real[keyy.shape[0] - pre
    Days + i] * 100))

110.
111.
112. # ----- #
113. # 作图
114. x1 = np.arange(1, keyy.shape[0] + 1, 1)
115. y1 = real # 实际数据
116. x2 = np.arange(keyy.shape[0] - preDays, keyy.shape[0] + 1, 1)
117. y2 = pre1 # 预测数据
118. x3 = np.arange(keyy.shape[0] - preDays, keyy.shape[0] + 1, 1)
119. y3 = pre2 # 预测数据
120. x4 = np.arange(keyy.shape[0] - preDays, keyy.shape[0] + 1, 1)
121. y4 = pre3 # 预测数据
122.
123. ax.plot(x1, y1, marker='o', linestyle='-', alpha=0.5, color='r') # 作图
124. ax.plot(x2, y2, marker='o', linestyle='--', alpha=0.5, color='b') # 作图
125. ax.plot(x3, y3, marker='o', linestyle='--', alpha=0.5, color='g') # 作图
126. ax.plot(x4, y4, marker='o', linestyle='--', color='y') # 作图
127.
128. xticks = [] # 构造 x 标签
129. for i in range(keyy.shape[0]):
130.     tmp1ist = keyy.iloc[i].tolist()
131.     xticks.append(str(tmp1ist[0])+'.'+str(tmp1ist[1])+'.'+str(tmp1ist[2]))
132.     pass
133. ax.set_xticks(range(1, keyy.shape[0] + 1)) # 设置 x 轴刻度
134. ax.set_xticklabels(xticks) # 设置 x 轴刻度标签
135. ax.set_xlabel("时间") # 设置 x 轴上的标签
136. ax.set_ylabel("人数") # 设置 y 轴上的标签
137. ax.set_title("全球变化趋势(预测{}天)".format(preDays), fontsize=12, color='grey') # 设置标题
138. first = True
139. for a, b in zip(x1, y1):
140.     ax.text(a - 0.25, b + 300000, '%d' % b, ha='center', va='center', fontsize=10, color='r') #
    数值点旁添加具体值
141. for a, b in zip(x2, y2):
142.     if first:
143.         first = False
144.     else:
145.         ax.text(a + 0.25, b - 300000, '%d' % b, ha='center', va='center', fontsize=10, color='b')
    # 数值点旁添加具体值
146. first = True
147. for a, b in zip(x3, y3):
148.     if first:

```

```
149.         first = False
150.     else:
151.         ax.text(a + 0.25, b - 200000, '%d' % b, ha='center', va='center', fontsize=10, color='g')
            # 数值点旁添加具体值
152. first = True
153. for a, b in zip(x4, y4):
154.     if first:
155.         first = False
156.     else:
157.         ax.text(a + 0.25, b - 400000, '%d' % b, ha='center', va='center', fontsize=10, color='y')
            # 数值点旁添加具体值
158. ax.legend(['实际', 'GM(1,1)预测', '差分预测', '一元线性回归'], loc='upper left') # 设置图例
159. plt.grid(True, linestyle='--', alpha=0.3, color='grey') # 设置网格线
160. plt.get_current_fig_manager().window.showMaximized() # 窗口最大化
161. plt.show()
```