

北京邮电大学

实验报告



题 目： 分析新冠疫情数据实验报告

课程名称：	Python程序设计
上课学期：	2021 春
授课教师：	杨亚
姓 名：	刘俊杰
学 号：	2019213687
日 期：	2021年12月31日

目录

- 1 实验内容
- 2 数据来源
- 3 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(`totCase.py`)
 - 3.2 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(`newCase.py`)
 - 3.2.1 15 天中每日新增确诊数累计排名前 10 个国家
 - 3.2.2 每日新增确诊数据的曲线图
 - 3.3 累计确诊数排名前 10 的国家名称及其数量(`totCaseTop10.py`)
 - 3.4 各个国家的累计确诊人数的比例(`totCastProportion.py`)
 - 3.5 累计确诊人数占国家总人口比例最高的 10 个国家(`caseRate.py`)
 - 3.6 疫苗接种情况 (`vaccineMap.py`)
 - 3.7 疫苗接种率 (累计疫苗接种/人数国家人数) 最低的 10 个国家(`vaccineRateLow.py`)
 - 3.8 全球 GDP 前十名国家的累计确诊人数箱型图(`totCaseInCDPTop10.py`)
 - 3.9 死亡率最高的 10 个国家(`fatalityRate.py`)
 - 3.10 其它你希望分析和展示的数据
 - 3.10.1 累计确诊人数占国家总人口比例最低的 10 个国家(`caseRateLow10.py`)
 - 3.10.2 疫苗接种率 (累计疫苗接种/人数国家人数) 最高的 10 个国家 (`vaccineRateTop.py`)
 - 3.10.3 死亡率最低的 10 个国家 (`fatalityRateLow.py`)
- 4 全世界应对新冠疫情最好的 10 个国家
- 5 预测(`totCasePred.py`)
- 6 源程序
 - 6.1 爬虫
 - 6.1.1 `spider.py`
 - 6.1.2 `items.py`
 - 6.1.3 `middlewares.py`
 - 6.1.4 `pipelines.py`
 - 6.1.5 `settings.py`
 - 6.2 数据处理及展示
 - 6.2.1 `totCase.py`
 - 6.2.2 `newCase.py`
 - 6.2.3 `totCaseTop10.py`
 - 6.2.4 `totCastProportion.py`
 - 6.2.5 `caseRate.py`
 - 6.2.6 `vaccineMap.py`
 - 6.2.7 `vaccineRateLow.py`
 - 6.2.8 `totCaseInCDPTop10.py`
 - 6.2.9 `fatalityRate.py`
 - 6.2.10 `caseRateLow10.py`
 - 6.2.11 `vaccineRateTop.py`
 - 6.2.12 `fatalityRateLow.py`
 - 6.2.13 `totCasePred.py`
- 7 总结

1 实验内容

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 2021 年 12 月 5 日开始的连续 15 天的数据，国家数不少于 100 个。

2 数据来源

新冠病毒数据网址 URL: <https://ourworldindata.org/coronavirus-data>

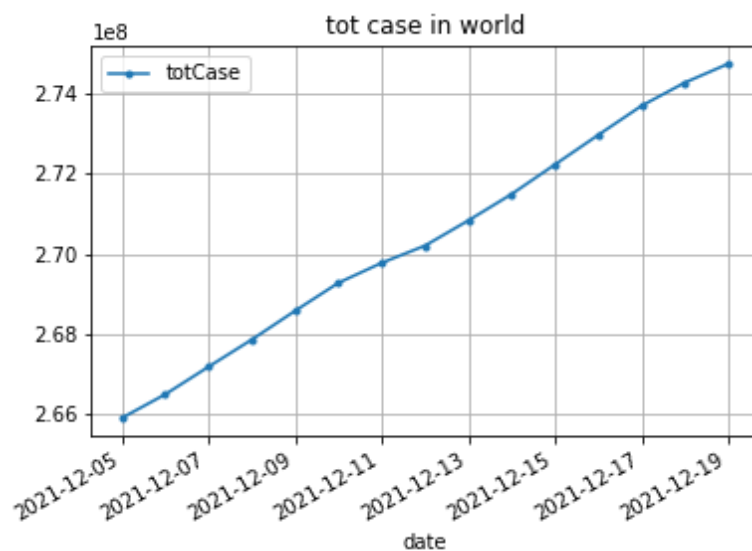
新冠疫苗数据网址 URL: https://ourworldindata.org/covid-vaccinations?country=OWID_WRL

首页截图:



3 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(totCase.py)



3.2 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(newCase.py)

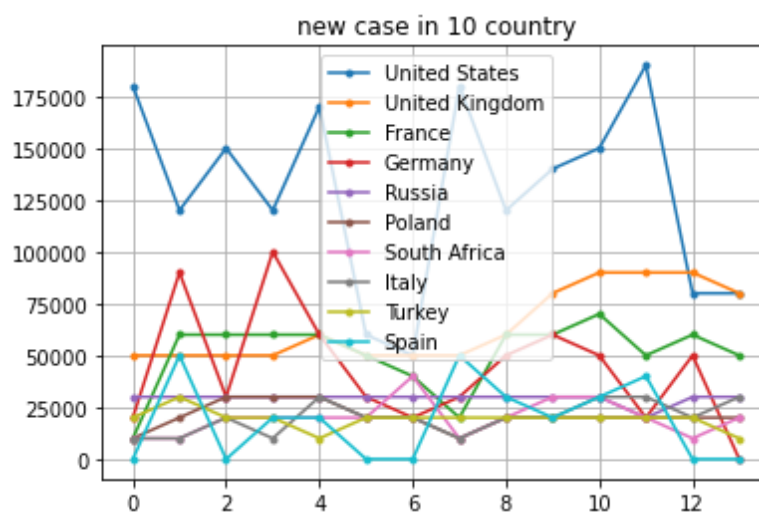
3.2.1 15 天中每日新增确诊数累计排名前 10 个国家

```

1 country,0
2 United States,1790000
3 United Kingdom,900000
4 France,710000
5 Germany,610000
6 Russia,410000
7 Poland,290000
8 South Africa,280000
9 Italy,280000
10 Turkey,270000
11 Spain,260000

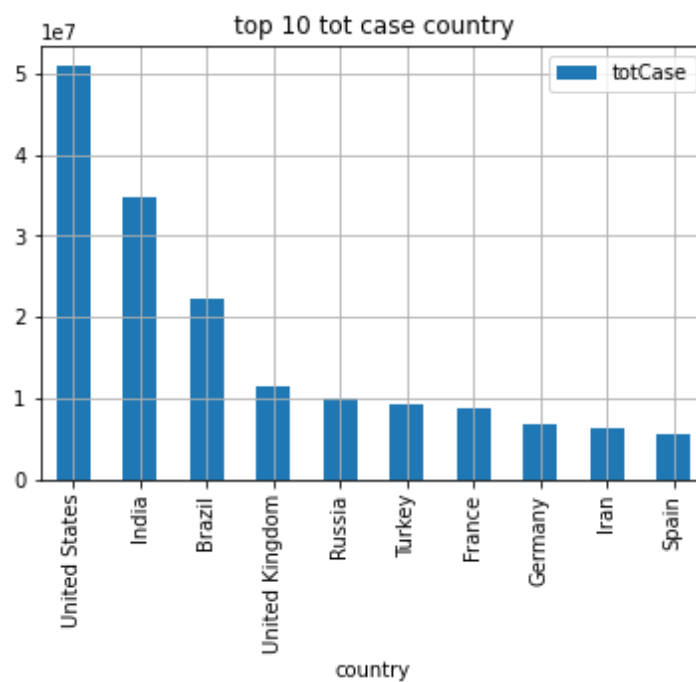
```

3.2.2 每日新增确诊数据的曲线图



3.3 累计确诊数排名前 10 的国家名称及其数量(totCaseTop10.py)

```
1 ,country,totCase,date
2 8,United States,50890000,2021-12-19
3 10,India,34750000,2021-12-19
4 11,Brazil,22220000,2021-12-19
5 12,United Kingdom,11380000,2021-12-19
6 13,Russia,10040000,2021-12-19
7 15,Turkey,9170000,2021-12-19
8 16,France,8670000,2021-12-19
9 17,Germany,6810000,2021-12-19
10 18,Iran,6170000,2021-12-19
11 19,Spain,5460000,2021-12-19
```

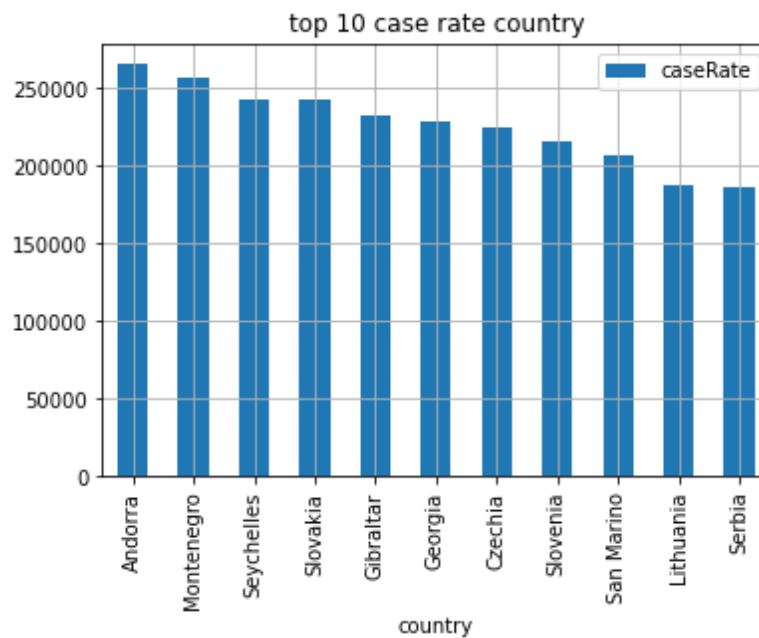


3.4 各个国家的累计确诊人数的比例(totCastProportion.py)

用饼图展示各个国家的累计确诊人数的比例（爬取的所有国家，数据较小的国家可以合并处理）



3.5 累计确诊人数占国家总人口比例最高的 10 个国家(caseRate.py)



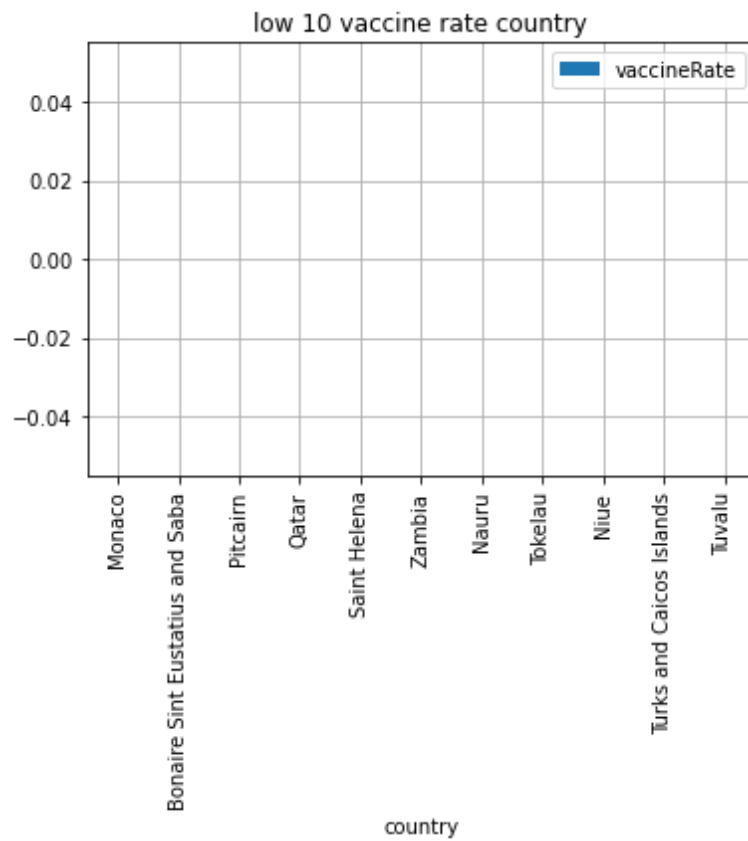
单位为每百万人确诊人数。

3.6 疫苗接种情况 (vaccineMap.py)

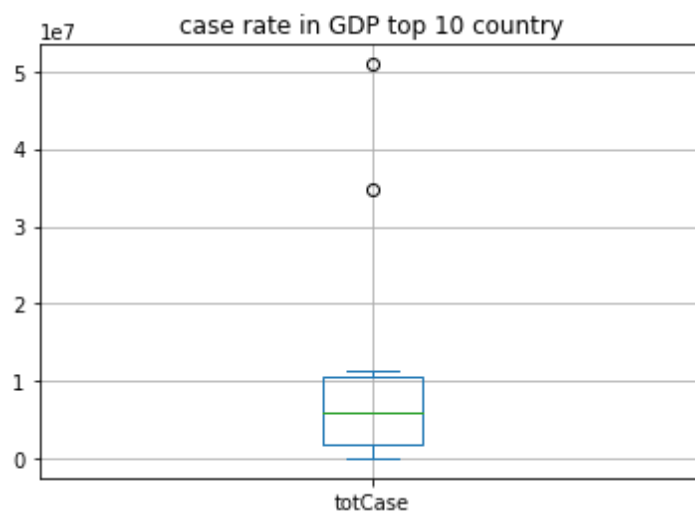
疫苗接种情况（至少接种了一针及以上），请用地图形式展示；



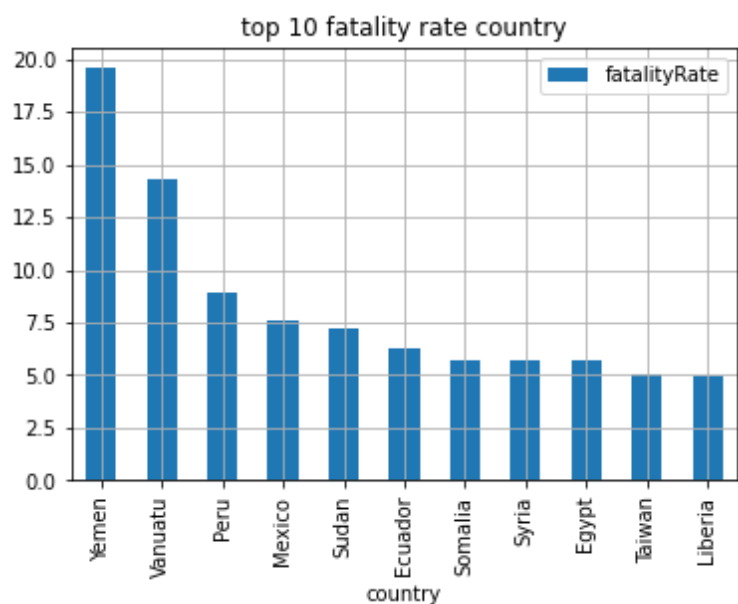
3.7 疫苗接种率（累计疫苗接种/人数国家人数）最低的 10 个国家 (vaccineRateLow.py)



3.8 全球 GDP 前十名国家的累计确诊人数箱型图(totCaseInCDPTop10.py)



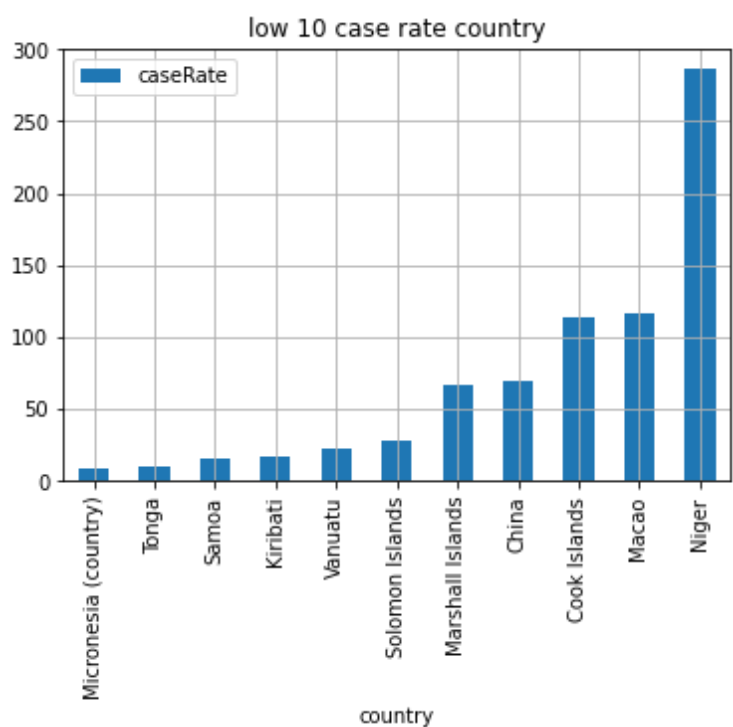
3.9 死亡率最高的 10 个国家(fatalityRate.py)



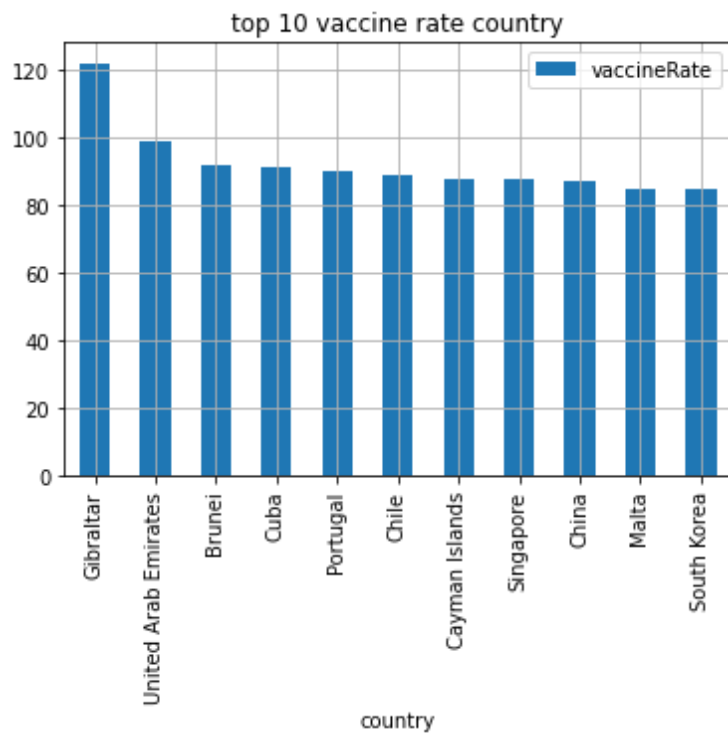
3.10 其它你希望分析和展示的数据

以上图形应包括完整的坐标、刻度、标签、图例等，如有必要请配上说明文字，对图中的内容进行解释。

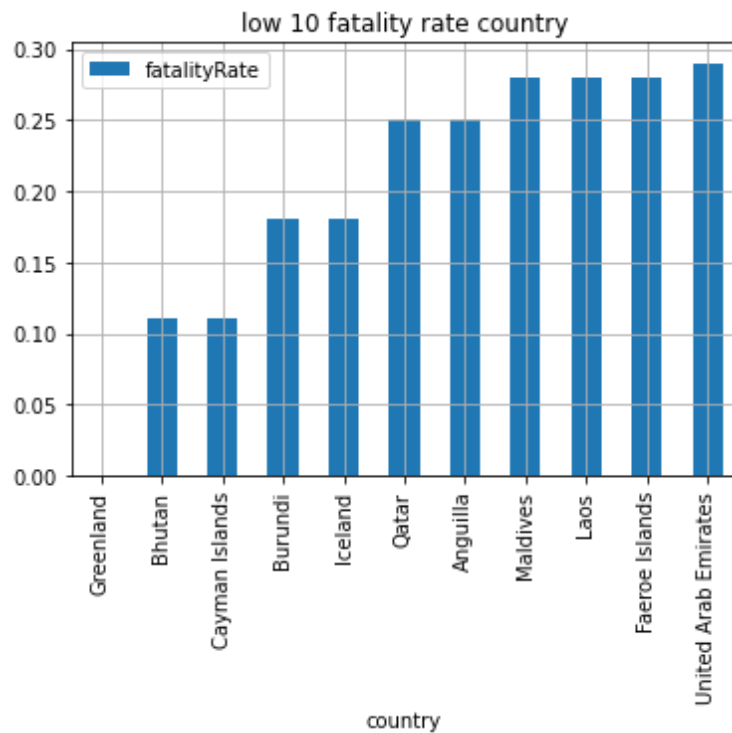
3.10.1 累计确诊人数占国家总人口比例最低的 10 个国家(caseRateLow10.py)



3.10.2 疫苗接种率（累计疫苗接种/人数国家人数）最高的 10 个国家（vaccineRateTop.py）



3.10.3 死亡率最低的 10 个国家 (fatalityRateLow.py)



4 全世界应对新冠疫情最好的 10 个国家

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。

中国确诊人数占国家总人口比例低，疫苗接种率（累计疫苗接种/人数国家人数）高，当之无愧对新冠疫情最好的 10 个国家之一。

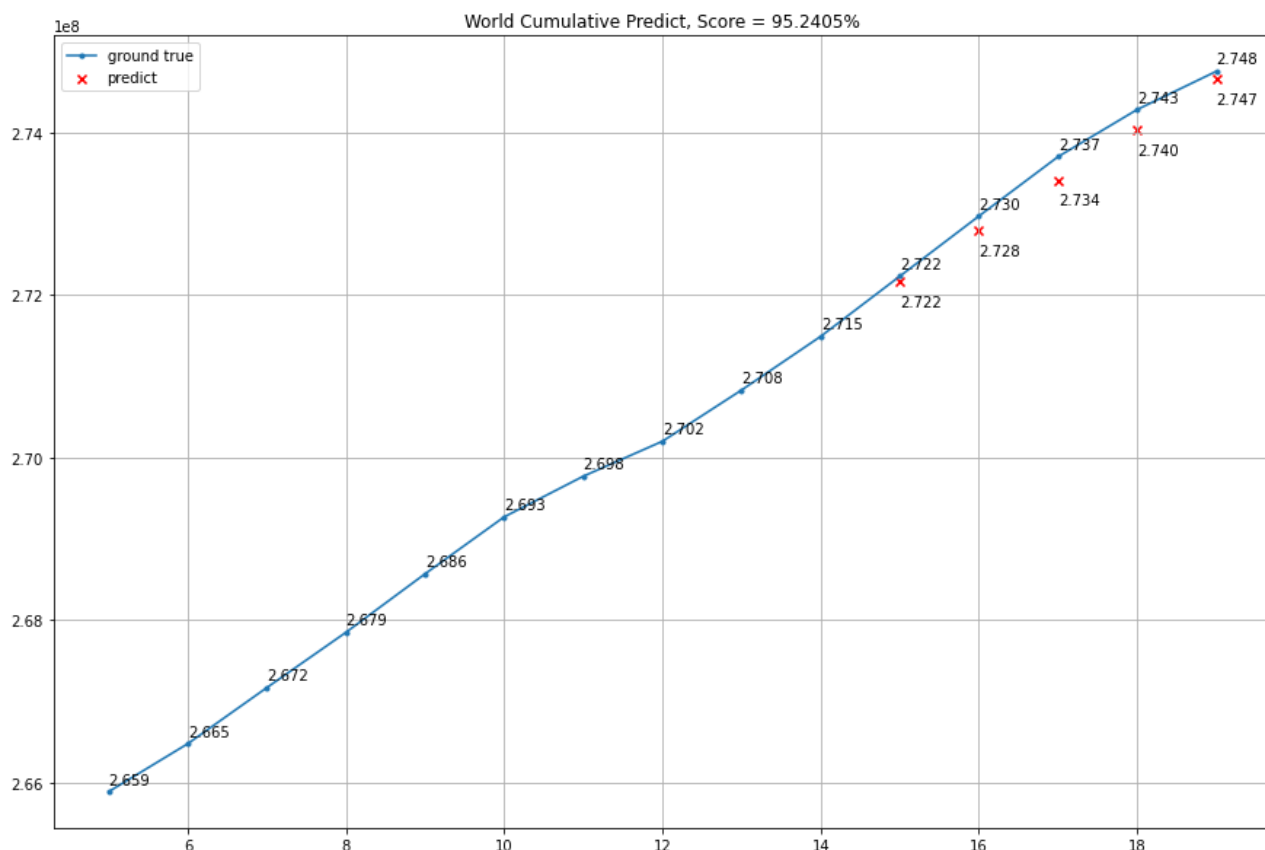
世界主要国家中还有墨西哥确诊率低，阿联酋、古巴、葡萄牙、新加坡、韩国等确诊率低，格陵兰岛、冰岛、老挝、阿联酋等死亡率低。

所以全世界应对新冠疫情最好的 10 个国家是：中国、墨西哥、阿联酋、古巴、葡萄牙、新加坡、韩国、格陵兰岛、冰岛、老挝。

5 预测(totCasePred.py)

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，并与实际数据进行对比。说明你预测的方法，并分析与实际数据的差距和原因。

使用 sklearn 中的 linear_model 线性回归模型对确诊数进行预测，真实数据为蓝色点，预测值为红色叉，可以发现预测值比真实值要小一些，产生差距的原因是线性回归模型是比较简单的模型，无法对如此复杂的全球累计确诊数进行精确的预测，现在的预测值已经非常接近真实值。



6 源程序

6.1 爬虫

6.1.1 spider.py

```
1 import scrapy
2 import re
3 from covid19.items import MyItem # 从 items.py 中引入 MyItem 对象
4
5 class mySpider(scrapy.spiders.Spider):
6     name = "covid19" # 爬虫的名字是 covid19
7     allow_domains = ['https://ourworldindata.org/explorers/'] # 允许爬取的网站
8     # totCase
9     start_urls = ['https://ourworldindata.org/explorers/']
```

```

10         'coronavirus-data-explorer?tab=table&zoomToSelection=true&time=2021-
12-' + str(date).zfill(2) +
11
12         '&facet=none&pickerSort=desc&pickerMetric=total_vaccinations_per_hundred&'\
13         'Metric=Confirmed+cases&Interval=Cumulative&Relative+to+Population=false&'\
14         'Color+by+test+positivity=false' for date in range(5, 20)] #爬取5~19日
数据
15     # caseRate
16     #start_urls = ['https://ourworldindata.org/explorers/'\
17     #     'coronavirus-data-explorer?tab=table&time=2021-12-19&facet=none&'\
18     #
19     'Metric=Confirmed+cases&Interval=Cumulative&Relative+to+Population=true&Color+by+test+positivity=false']
20
21     # vaccineRate
22     #start_urls = ['https://ourworldindata.org/explorers/'\
23     #     'coronavirus-data-explorer?tab=table&time=2021-12-19&facet=none&'\
24     #
25     'Metric=People+vaccinated&Interval=Cumulative&Relative+to+Population=true&Color+by+test+positivity=false']
26
27     # fatalityRate
28     #start_urls = ['https://ourworldindata.org/explorers/'\
29     #     'coronavirus-data-explorer?tab=table&time=2021-12-19&facet=none&'\
30     #
31     'uniformYAxis=0&pickerSort=asc&pickerMetric=location&Metric=Case+fatality+rate&Interval=Cumulative&Relative+to+Population=true&Color+by+test+positivity=false&country=~OWID_WRL']
32
33     def parse(self, response):
34         item = MyItem()
35         for each in
response.xpath('/html/body/main/div/div[3]/div/div[1]/div/table/tbody/tr'):
36             item['country'] = each.xpath('td[1]/text()').get()
37             # totCase
38
39             item['totCase'] = each.xpath('td[2]/text()').get().replace(',', '',
'') #去除逗号
40             if re.search('million', item['totCase']): #将million转换为数字概念的
1e6
41                 item['totCase'] = item['totCase'].replace(' million', '')
42                 item['totCase'] = int(float(item['totCase']) * 1000000)
43             else:
44                 item['totCase'] = int(item['totCase'])
45
46         # caseRate

```

```

44         # item['caseRate'] =
float(each.xpath('td[2]/text()').get().replace(',','')) #去除逗号
45         # vaccineRate
46         #item['vaccineRate'] = each.xpath('td[2]/text()').get()
47         #if item['vaccineRate'] == None: # 没有接种疫苗的国家或地区不考虑
48         #     continue
49         #item['vaccineRate'] = item['vaccineRate'].replace(' ','')
50         # fatalityRate
51         #item['fatalityRate'] = each.xpath('td[2]/text()').get()
52         #if item['fatalityRate'] == None: #没有死亡率的国家认为是0%
53         #     item['fatalityRate'] = '0.00%'
54
55         item['date'] = re.findall(".*time=(.*?)&.*", response.url)[0] #获
取爬取的日期
56         yield item

```

6.1.2 items.py

```

1 # Define here the models for your scraped items
2 #
3 # See documentation in:
4 # https://docs.scrapy.org/en/latest/topics/items.html
5
6 import scrapy
7
8 class MyItem(scrapy.Item):
9     # define the fields for your item here like:
10     # name = scrapy.Field()
11     country = scrapy.Field() #国家或地区名
12     date = scrapy.Field() #日期
13     # totCase
14     totCase = scrapy.Field() #累计确诊病例
15     # caseRate
16     #caseRate = scrapy.Field() #确诊病例占总人口比
17     # vaccineRate
18     #vaccineRate = scrapy.Field() #接种至少一针疫苗人数占总人口比
19     # fatalityRate
20     #fatalityRate = scrapy.Field() #死亡率
21     pass

```

6.1.3 middlewares.py

```

1 # Define here the models for your spider middleware
2 #
3 # See documentation in:
4 # https://docs.scrapy.org/en/latest/topics/spider-middleware.html
5
6 from scrapy import signals

```

```

7
8 # useful for handling different item types with a single interface
9 from itemadapter import is_item, ItemAdapter
10 from selenium import webdriver
11 from scrapy.http import HtmlResponse
12 from time import sleep
13
14 class Covid19SpiderMiddleware:
15     # Not all methods need to be defined. If a method is not defined,
16     # scrapy acts as if the spider middleware does not modify the
17     # passed objects.
18
19     @classmethod
20     def from_crawler(cls, crawler):
21         # This method is used by Scrapy to create your spiders.
22         s = cls()
23         crawler.signals.connect(s.spider_opened,
signal=signals.spider_opened)
24         return s
25
26     def process_spider_input(self, response, spider):
27         # Called for each response that goes through the spider
28         # middleware and into the spider.
29
30         # Should return None or raise an exception.
31         return None
32
33     def process_spider_output(self, response, result, spider):
34         # Called with the results returned from the Spider, after
35         # it has processed the response.
36
37         # Must return an iterable of Request, or item objects.
38         for i in result:
39             yield i
40
41     def process_spider_exception(self, response, exception, spider):
42         # Called when a spider or process_spider_input() method
43         # (from other spider middleware) raises an exception.
44
45         # Should return either None or an iterable of Request or item
objects.
46         pass
47
48     def process_start_requests(self, start_requests, spider):
49         # Called with the start requests of the spider, and works
50         # similarly to the process_spider_output() method, except
51         # that it doesn't have a response associated.
52

```

```

53         # Must return only requests (not items).
54         for r in start_requests:
55             yield r
56
57     def spider_opened(self, spider):
58         spider.logger.info('Spider opened: %s' % spider.name)
59
60
61 class Covid19DownloaderMiddleware:
62     # Not all methods need to be defined. If a method is not defined,
63     # scrapy acts as if the downloader middleware does not modify the
64     # passed objects.
65
66     @classmethod
67     def from_crawler(cls, crawler):
68         # This method is used by Scrapy to create your spiders.
69         s = cls()
70         crawler.signals.connect(s.spider_opened,
signal=signals.spider_opened)
71         return s
72
73     def process_request(self, request, spider):
74         # Called for each request that goes through the downloader
75         # middleware.
76
77         # Must either:
78         # - return None: continue processing this request
79         # - or return a Response object
80         # - or return a Request object
81         # - or raise IgnoreRequest: process_exception() methods of
82         #   installed downloader middleware will be called
83
84         options = webdriver.ChromeOptions()
85         options.add_argument('--headless')
86         options.add_argument("--ignore-certificate-error")
87         options.add_argument("--ignore-ssl-errors")
88         driver = webdriver.Chrome(chrome_options=options)
89         driver.get(request.url)
90         driver.implicitly_wait(120)
91         sleep(10)
92         html = driver.page_source
93         return HtmlResponse(url=driver.current_url, body=html,
request=request, encoding='utf-8')
94
95     def process_response(self, request, response, spider):
96         # Called with the response returned from the downloader.
97
98         # Must either;

```

```

99         # - return a Response object
100        # - return a Request object
101        # - or raise IgnoreRequest
102        return response
103
104    def process_exception(self, request, exception, spider):
105        # Called when a download handler or a process_request()
106        # (from other downloader middleware) raises an exception.
107
108        # Must either:
109        # - return None: continue processing this exception
110        # - return a Response object: stops process_exception() chain
111        # - return a Request object: stops process_exception() chain
112        pass
113
114    def spider_opened(self, spider):
115        spider.logger.info('Spider opened: %s' % spider.name)
116

```

6.1.4 pipelines.py

```

1  from itemadapter import ItemAdapter
2  import json
3
4  class MyPipeline:
5      def open_spider(self, spider):
6          try:
7              self.file = open('totCase.json', 'w', encoding='utf-8')
8              #self.file = open('caseRate.json', 'w', encoding='utf-8')
9              #self.file = open('vaccineRate.json', 'w', encoding='utf-8')
10             #self.file = open('fatalityRate.json', 'w', encoding='utf-8')
11             self.file.write('[')
12         except Exception as err:
13             print(err)
14
15     def process_item(self, item, spider):
16         dict_item = dict(item) #生成字典对象
17         json_str = json.dumps(dict_item, ensure_ascii=False) + ",\n" #生成
json串
18         self.file.write(json_str) #将json串写入到文件中
19         return item
20
21     def close_spider(self, spider):
22         self.file.write(']')
23         self.file.close() #关闭文件

```


6.1.5 settings.py

```
1 # Scrapy settings for covid19 project
2 #
3 # For simplicity, this file contains only settings considered important or
4 # commonly used. You can find more settings consulting the documentation:
5 #
6 #     https://docs.scrapy.org/en/latest/topics/settings.html
7 #     https://docs.scrapy.org/en/latest/topics/downloader-middleware.html
8 #     https://docs.scrapy.org/en/latest/topics/spider-middleware.html
9
10 BOT_NAME = 'covid19'
11
12 SPIDER_MODULES = ['covid19.spiders']
13 NEWSPIDER_MODULE = 'covid19.spiders'
14
15
16 # Crawl responsibly by identifying yourself (and your website) on the user-
17 # agent
18 #USER_AGENT = 'covid19 (+http://www.yourdomain.com)'
19
20 # Obey robots.txt rules
21 # ROBOTSTXT_OBEY = True
22 ROBOTSTXT_OBEY = False
23
24 # Configure maximum concurrent requests performed by Scrapy (default: 16)
25 #CONCURRENT_REQUESTS = 32
26
27 # Configure a delay for requests for the same website (default: 0)
28 # See https://docs.scrapy.org/en/latest/topics/settings.html#download-delay
29 # See also autothrottle settings and docs
30 DOWNLOAD_DELAY = 3
31 # The download delay setting will honor only one of:
32 #CONCURRENT_REQUESTS_PER_DOMAIN = 16
33 #CONCURRENT_REQUESTS_PER_IP = 16
34
35 # Disable cookies (enabled by default)
36 #COOKIES_ENABLED = False
37
38 # Disable Telnet Console (enabled by default)
39 #TELNETCONSOLE_ENABLED = False
40
41 # Override the default request headers:
42 #DEFAULT_REQUEST_HEADERS = {
43 #    'Accept':
44 #        'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
45 #    'Accept-Language': 'en',
46 #}
```

```
45
46 # Enable or disable spider middlewares
47 # See https://docs.scrapy.org/en/latest/topics/spider-middleware.html
48 # SPIDER_MIDDLEWARES = {
49 #     'covid19.middlewares.Covid19SpiderMiddleware': 543,
50 # }
51
52 # Enable or disable downloader middlewares
53 # See https://docs.scrapy.org/en/latest/topics/downloader-middleware.html
54 DOWNLOADER_MIDDLEWARES = {
55     'covid19.middlewares.Covid19DownloaderMiddleware': 543,
56 }
57
58 # Enable or disable extensions
59 # See https://docs.scrapy.org/en/latest/topics/extensions.html
60 #EXTENSIONS = {
61 #     'scrapy.extensions.telnet.TelnetConsole': None,
62 #}
63
64 # Configure item pipelines
65 # See https://docs.scrapy.org/en/latest/topics/item-pipeline.html
66 #ITEM_PIPELINES = {
67 #     'covid19.pipelines.Covid19Pipeline': 300,
68 #}
69 ITEM_PIPELINES = {
70     'covid19.pipelines.MyPipeline': 300,
71 }
72
73 # Enable and configure the AutoThrottle extension (disabled by default)
74 # See https://docs.scrapy.org/en/latest/topics/autothrottle.html
75 AUTOTHROTTLER_ENABLED = True
76 # The initial download delay
77 AUTOTHROTTLER_START_DELAY = 5
78 # The maximum download delay to be set in case of high latencies
79 AUTOTHROTTLER_MAX_DELAY = 60
80 # The average number of requests Scrapy should be sending in parallel to
81 # each remote server
82 AUTOTHROTTLER_TARGET_CONCURRENCY = 1.0
83 # Enable showing throttling stats for every response received:
84 #AUTOTHROTTLER_DEBUG = False
85
86 # Enable and configure HTTP caching (disabled by default)
87 # See https://docs.scrapy.org/en/latest/topics/downloader-
88 # middleware.html#httpcache-middleware-settings
89 #HTTPCACHE_ENABLED = True
90 #HTTPCACHE_EXPIRATION_SECS = 0
91 #HTTPCACHE_DIR = 'httpcache'
92 #HTTPCACHE_IGNORE_HTTP_CODES = []
```

```

92 #HTTPCACHE_STORAGE = 'scrapy.extensions.httpcache.FilesystemCacheStorage'
93

```

6.2 数据处理及展示

6.2.1 totCase.py

```

1  #%%
2  import pandas as pd
3  totCase = pd.read_json('./covid19/totCase.json')
4  totCaseWrold = totCase.loc[totCase['country'] == 'World'].drop(columns=
    ['country']).set_index('date')
5  totCaseWrold.plot(marker='.', grid=True)
6  # %%
7

```

6.2.2 newCase.py

```

1  #%%
2  import pandas as pd
3  totCase = pd.read_json('./covid19/totCase.json')
4  newCase = totCase.groupby('country')['totCase'].agg(['min', 'max'])
5  newCase = (newCase['max'] - newCase['min']).sort_values(ascending=False)
6  print(newCase)
7  newCase.to_csv('newCase.csv', index=True)
8  newCases = []
9  countrys = ['United States', 'United
    Kingdom', 'France', 'Germany', 'Russia', 'Poland', 'South
    Africa', 'Italy', 'Turkey', 'Spain']
10 for country in countrys:
11     newCase = totCase.loc[totCase['country'] == country]
12     newCase = newCase.drop(columns=['country']).set_index('date')
13     newCase = newCase.loc[['2021-12-' + str(date).zfill(2) for date in
    range(6, 20)]] .reset_index(drop=True) . \
14         subtract(newCase.loc[['2021-12-' + str(date).zfill(2) for date in
    range(5, 19)]] .reset_index(drop=True))
15     newCase = newCase.rename({'totCase':country}, axis=1)
16     newCases.append(newCase)
17 df = pd.concat(newCases)
18 df.plot(title='new case in 10 country', marker='.', grid=True)
19 # %%
20

```

6.2.3 totCaseTop10.py

```

1  ###
2  import pandas as pd
3  totCase = pd.read_json('./covid19/totCase.json')
4  totCaseTop = totCase.loc[totCase['date'] == '2021-12-19'].sort_values(by=
    ['totCase'], ascending=False).reset_index(drop=True)
5  print(totCaseTop)
6  totCaseTop.to_csv('totCaseTop10.csv', index=True)
7  Top10Country = ["United States", "India", "Brazil", "United
    Kingdom", "Russia", "Turkey", "France", "Germany", "Iran", "Spain"]
8  totCaseTop10 = totCaseTop.loc[totCaseTop['country'].isin(Top10Country)]
9  totCaseTop10 = totCaseTop10.set_index('country')[['totCase']]
10 print(totCaseTop10)
11 totCaseTop10.plot.bar(title='top 10 tot case country', grid=True)
12 # %%
13

```

6.2.4 totCastProportion.py

```

1  ###
2  import pandas as pd
3  totCase = pd.read_json('./covid19/totCase.json')
4  totCaseTop = totCase.loc[totCase['date'] == '2021-12-19'].sort_values(by=
    ['totCase'], ascending=False).reset_index(drop=True).set_index('country')
    [['totCase']]
5  dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
    'Lower middle income', 'North America', 'European Union', 'South America',
    'Africa']
6  totCaseTop = totCaseTop.drop(dropList, axis=0).reset_index()
7  totCaseCount = totCaseTop.loc[0:10,:]
8  totCaseCount = totCaseCount.set_index('country')
9  totCaseCount.loc['other'] = totCaseTop.loc[10:].sum(axis=1)
10 totCaseCount.plot.pie(title='tot case pie in country', y='totCase', figsize=
    (5, 5))
11 # %%
12

```

6.2.5 caseRate.py

```

1  ###
2  import pandas as pd
3  caseRate = pd.read_json('./covid19/caseRate.json')
4  caseRateTop = caseRate.sort_values(by=['caseRate'],
   ascending=False).set_index('country')
5  print(caseRateTop)
6  dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
   'Lower middle income', 'North America', 'European Union', 'South America',
   'Africa']
7  caseRateTop = caseRateTop.drop(dropList,
   axis=0).reset_index().loc[0:10,:].set_index('country')[['caseRate']]
8  caseRateTop.plot.bar(title='top 10 case rate country', grid=True)
9  ###
10

```

6.2.6 vaccineMap.py

```

1  ###
2  import pandas as pd
3  from pyecharts import options as opts
4  from pyecharts.charts import Map
5  import imgkit
6  vaccineRate = pd.read_json('./covid19/vaccineRate.json')
7  print(vaccineRate)
8  for i in vaccineRate.index:
9      vaccineRate['vaccineRate'][i] = float(vaccineRate['vaccineRate']
   [i].replace('%', ''))
10 vaccineRate = vaccineRate.sort_values(by=['vaccineRate'],
   ascending=False).set_index('country')
11 dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
   'Lower middle income', 'North America', 'European Union', 'South America',
   'Africa']
12 vaccineRate = vaccineRate.drop(dropList, axis=0)
   [['vaccineRate']].reset_index()
13
14 countries = list([vaccineRate['country'][i] for i in vaccineRate.index])
15 vaccineRates = list([vaccineRate['vaccineRate'][i] for i in
   vaccineRate.index])
16 map = Map(init_opts=opts.InitOpts(width='100%', height='1200px'))
17 map.add("vaccineRate", [list(z) for z in zip(countries, vaccineRates)],
   'world')
18 map.set_global_opts(
19     title_opts=opts.TitleOpts(title='World VaccineRate'), #标题
20     visualmap_opts=opts.VisualMapOpts(min_=0,max_=100)) #热力图数值区间
21 map.set_series_opts(label_opts=opts.LabelOpts(is_show=True)) #热力图图例
22 map.render("map.html") #导出html
23 imgkit.from_file('./map.html', './map.jpg')

```

```
24 # %%
25
```

6.2.7 vaccineRateLow.py

```
1  ###
2  import pandas as pd
3  vaccineRate = pd.read_json('./covid19/vaccineRate.json')
4  print(vaccineRate)
5  for i in vaccineRate.index:
6      vaccineRate['vaccineRate'][i] = float(vaccineRate['vaccineRate']
7      [i].replace('%', ''))
8  vaccineRateLow = vaccineRate.sort_values(by=
9  ['vaccineRate']).set_index('country')
10 dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
11 'Lower middle income', 'North America', 'European Union', 'South America',
12 'Africa']
13 vaccineRateLow = vaccineRateLow.drop(dropList,
14 axis=0).reset_index().loc[0:10,:].set_index('country')[['vaccineRate']]
15 vaccineRateLow.plot.bar(title='low 10 vaccine rate country', grid=True)
16 # %%
17
```

6.2.8 totCaseInCDPTop10.py

```
1  ###
2  import pandas as pd
3  caseRate = pd.read_json('./covid19/totCase.json')
4  GDPTop10 = ['United States', 'China', 'Japan', 'Germany', 'United Kingdom',
5  'India', 'France', 'Italy', 'Canada', 'South Korea']
6  caseRate = caseRate.loc[caseRate['date'] == '2021-12-19']
7  caseRate = caseRate.loc[caseRate['country'].isin(GDPTop10)]
8  print(caseRate)
9  caseRate = caseRate.set_index('country')['totCase']
10 caseRate.plot.box(title='case rate in GDP top 10 country', grid=True)
11 # %%
12
```

6.2.9 fatalityRate.py

```

1  """
2  import pandas as pd
3  fatalityRate = pd.read_json('./covid19/fatalityRate.json')
4  for i in fatalityRate.index:
5      fatalityRate['fatalityRate'][i] = float(fatalityRate['fatalityRate']
6      [i].replace('%', ''))
7  fatalityRateTop = fatalityRate.sort_values(by=['fatalityRate'],
8  ascending=False).reset_index(drop=True).set_index('country')
9  dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
10 'Lower middle income', 'North America', 'European Union', 'South America',
11 'Africa']
12 fatalityRateTop = fatalityRateTop.drop(dropList,
13 axis=0).reset_index().loc[0:10,:].set_index('country')[['fatalityRate']]
14 fatalityRateTop.plot.bar(title='top 10 fatality rate country', grid=True)
15 # """
16

```

6.2.10 caseRateLow10.py

```

1  """
2  import pandas as pd
3  caseRate = pd.read_json('./covid19/caseRate.json')
4  caseRateLow = caseRate.sort_values(by=['caseRate']).set_index('country')
5  dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
6 'Lower middle income', 'North America', 'European Union', 'South America',
7 'Africa']
8 caseRateLow = caseRateLow.drop(dropList,
9 axis=0).reset_index().loc[0:10,:].set_index('country')[['caseRate']]
10 caseRateLow.plot.bar(title='low 10 case rate country', grid=True)
11 # """
12

```

6.2.11 vaccineRateTop.py

```

1  ###
2  import pandas as pd
3  vaccineRate = pd.read_json('./covid19/vaccineRate.json')
4  print(vaccineRate)
5  for i in vaccineRate.index:
6      vaccineRate['vaccineRate'][i] = float(vaccineRate['vaccineRate']
7      [i].replace('%', ''))
8  vaccineRateLow = vaccineRate.sort_values(by=['vaccineRate'],
9  ascending=False).set_index('country')
10 dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
11 'Lower middle income', 'North America', 'European Union', 'South America',
12 'Africa']
13 vaccineRateLow = vaccineRateLow.drop(dropList,
14 axis=0).reset_index().loc[0:10,:].set_index('country')[['vaccineRate']]
15 vaccineRateLow.plot.bar(title='top 10 vaccine rate country', grid=True)
16 # %%
17

```

6.2.12 fatalityRateLow.py

```

1  ###
2  import pandas as pd
3  fatalityRate = pd.read_json('./covid19/fatalityRate.json')
4  for i in fatalityRate.index:
5      fatalityRate['fatalityRate'][i] = float(fatalityRate['fatalityRate']
6      [i].replace('%', ''))
7  fatalityRateTop = fatalityRate.sort_values(by=
8  ['fatalityRate']).reset_index(drop=True).set_index('country')
9  dropList = ['World', 'High income', 'Upper middle income', 'Asia', 'Europe',
10 'Lower middle income', 'North America', 'European Union', 'South America',
11 'Africa']
12 fatalityRateTop = fatalityRateTop.drop(dropList,
13 axis=0).reset_index().loc[0:10,:].set_index('country')[['fatalityRate']]
14 fatalityRateTop.plot.bar(title='low 10 fatality rate country', grid=True)
15 # %%
16

```

6.2.13 totCasePred.py

```

1  ###
2  import pandas as pd
3  import numpy as np
4  from sklearn import linear_model
5  from matplotlib import pyplot as plt
6  totCase = pd.read_json('./covid19/totCase.json')
7  for i in totCase.index:
8      totCase['date'][i] = int(str(totCase['date'][i]).split(' ')[0].split('-
9      ')[2])
10

```



```

9 totCaseWrold = totCase.loc[totCase['country'] == 'World'].drop(columns=
  ['country']).sort_values('date').set_index('date')
10 x = []
11 y = []
12 plt.figure(figsize=(15, 10))
13 for i in totCaseWrold.index:
14     x.append(i)
15     y.append(totCaseWrold['totCase'][i])
16 date = x.copy()
17 x2 = x[10:]
18 x = x[:10]
19 y2 = y[10:]
20 y = y[:10]
21 plt.grid(True)
22 plt.plot(date, y + y2, marker='.', label='ground true')
23 for a,b in zip(date, y + y2):
24     plt.text(a, b+1e5, '%.3f'%(b/1e8))
25 x2 = [_ for _ in range(15, 20)]
26 model = linear_model.LinearRegression()
27 model.fit(np.array(x).reshape(-1, 1), y)
28 y3 = model.predict(np.array(x2).reshape(-1, 1))
29 plt.scatter(date[10:], y3, c='red', marker='x', label='predict')
30 for a,b in zip(date[10:], y3):
31     plt.text(a, b-3e5, '%.3f'%(b/1e8))
32 plt.legend(loc=2)
33 score = model.score(np.array(x2).reshape(-1, 1), y2)
34 plt.title('World Cumulative Predict, Score = ' + str('%0.4f'%(score*100)) +
  '%')
35 plt.show()
36
37 #totCaseWrold.plot(title='tot case in world', marker='.', grid=True)
38 # %%
39

```

7 总结

通过爬虫得到新冠病毒和疫苗的数据并保存到 json 文件，使用 Python 的 Pandas 库对数据进行清洗、提取、排序和统计，然后使用 Dataframe 结构 plot 方法对数据进行可视化。要进行更进一步的可视化，可以使用 pyecharts 库绘制热力图。要利用现有的数据对未来进行预测，可以使用 Sklearn 库进行预测，该库包含了众多经典机器学习算法，可以很方便使用。

通过上述的工作，利用网络数据，可以让我们进一步分析和了解新冠疫情的情况，也能通过现有数据预测疫情的发展情况，有利于我们更好地应对疫情。通过这次实验，进一步熟悉了数据获取和分析的工具和方法，进一步加深了对 Python 的理解，受益匪浅。