

Python大作业：疫情数据分析

2019213678 谭海川 2019211301班

目录：



一、作业概述

二、数据爬取

1.数据来源

2.爬虫实现

分析网址

爬虫结构

爬虫核心代码spider.py

三、数据分析

1.用到的库

2.图表工具pyecharts

3.数据读入

4.按要求进行数据分析

四、实验总结

一、作业概述

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。

1. 标明你的数据来源：包括网址和首页截图

2. 数据分析和展示应包括：

1) 15 天中，全球新冠疫情的总体变化趋势；

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的曲线图；

3) 累计确诊数排名前 10 的国家名称及其数量；

4) 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）；

- 5) 累计确诊人数占国家总人口比例最高的 10 个国家；
- 6) 疫苗接种情况（至少接种了一针及以上），请用地图形式展示；
- 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；
- 8) 全球 GDP 前十名国家的累计确诊人数箱型图，要有平均值；
- 9) 死亡率最高的 10 个国家；
- 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. 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。

4. 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。

二、数据爬取

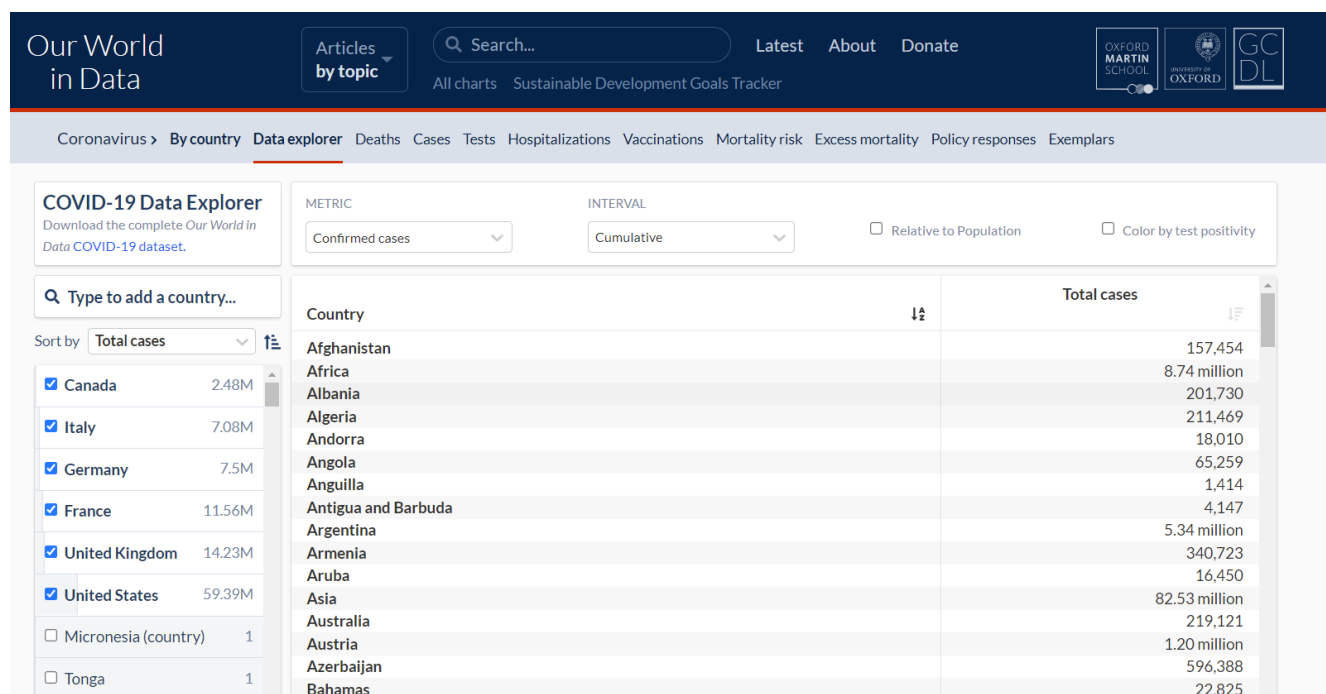
1. 数据来源

来自网站：点击打开：

COVID-19 Data Explorer

Research and data to make progress against the world's largest problems

<https://ourworldindata.org/explorers/coronavirus-data-explorer?tab=table&zoomToSelection=true&time=2021-12-05&f...>



该网站有按日计算的各个国家的：

- 确诊人数
- 新增确诊人数
- 死亡病例数
- 接种疫苗人数

数据齐全，可以通过之前学习的爬虫技术进行爬取

2.爬虫实现

分析网址

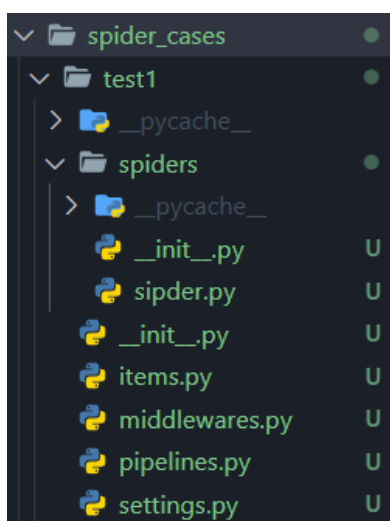
首先分析网页和网址结构，想办法找到历史数据。

观察某天的网址发现：

```
https://ourworldindata.org/explorers/coronavirus-data-explorer?
tab=table&zoomToSelection=true&time=2021-12-
05&facet=none&uniformYAxis=0&pickerSort=asc&pickerMetric=total_cases&Metric=Confirmed+cases&Interval=Cumulative&Relative+to+Population=false&Color+by+test+positivity=false
```

- Time字段可以输入日期，而metric字段可以选择网页展示的是确诊病例数/死亡数/疫苗接种数
- 由此，可以对于2021-12-05到2021-12-19的连续15天分别构造相应的url，按照这些url进行爬取。

爬虫结构



spider.py是核心代码，其余的部分和之前提交的爬虫代码大同小异，篇幅限制不再附代码。

爬虫核心代码spider.py

- 由于该网站还统计了诸如 世界 高收入国家 低收入国家 亚洲 欧洲等非国家级别的数据，所以我在爬取的时候进行判断，去掉这些数据。
- 对于这15天，会自动填充相应的url进行爬取

- 由于网站中数据超过1million和1billion的数据是用的million和billion表示，所以在爬取的时候分别乘1e6和1e9进行转化。但是这样导致的结果之一是很多国家统计的人数精度不够，因此，后续数据分析过程多数时候单位采用的是 百万人

Python

```
1 import scrapy
2 from scrapy.http import response
3 from test1.items import Test1Item
4 not_country=["World","High income","Upper middle income","Asia","Europe","Lower
   middle income","North America","European Union","South America","Africa","Ocean
   ia"]
5 class mySpider(scrapy.spiders.Spider):
6     name="bupt"
7     allowed_domains=["ourworldindata.org"]
8     with open("./date.txt","r") as f:
9         date=f.read()
10        #对于不同的15天, 会自动修改对应的网址进行爬取
11    start_urls=["https://ourworldindata.org/explorers/coronavirus-data-explorer?
   tab=table&zoomToSelection=true&time=2021-12-{}&facet=none&uniformYAxis=0&pickerS
   ort=asc&pickerMetric=total_cases&Metric=Confirmed+cases&Interval=Cumulative&Rela
   tive+to+Population=false&Color+by+test+positivity=false".format(date)]
12    def start_requests(self):
13        for url in self.start_urls:
14            yield scrapy.Request(url,callback=self.parse,cb_kwargs={"page_num":1
   })
15    def parse(self,response,**kwargs):
16        item=Test1Item()
17        for each in response.xpath("/html/body/main/div/div[3]/div/div[1]/div/ta
   ble/tbody/*"):
18            item['country']=each.xpath("./td[1]/text()").get()
19            cases=str(each.xpath("./td[2]/text()").get())
20            cases=cases.replace(",","")
21            #以下是对于输入为million的处理
22            if (cases.find("million")!=-1):
23                cases=cases.replace("million","")
24                item['total_cases']=int((float(cases)*1000000))
25            else :
26                item['total_cases']=int(cases)
27            if (item['country']and item['country'] not in not_country):
28                yield(item)
```

三、数据分析

1.用到的库

`csv` 库是用来进行csv文件的读入和输出

`scipy` 库是用来进行科学计算的库，我用来计算线性回归进行预测。

`pyecharts` 库是用来进行画图的库

Python

```
1 import pyecharts.options as opts
2 from pyecharts.charts import Line
3 from pyecharts.charts import Pie
4 from pyecharts.charts import Map
5 from pyecharts.charts import Boxplot
6 import scipy.stats as st
7 import csv
```

2.图表工具pyecharts

pyecharts 是一个用于生成 Echarts 图表的类库。Echarts 是百度开源的一个数据可视化 JS 库。用 Echarts 生成的图可视化效果非常棒，pyecharts 与 Python 进行对接，方便在 Python 中直接使用数据生成图。

3.数据读入

- cases是存储每个国家每天病例总数的字典，从爬取的csv文件中读取即可
- 死亡病例数，疫苗接种人数等读入的方法和总病例数是一样的。

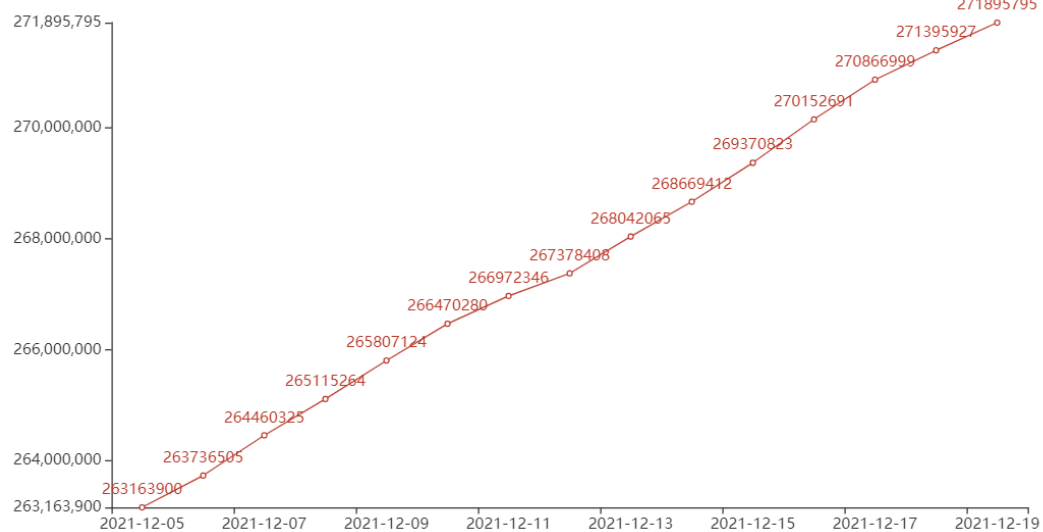
```
1 cases={}
2 country_list=[]
3 for date in date_list:
4     with open("./spider_cases/"+date+".csv","r") as f:
5         data=f.readlines()[1:]
6         for line in data:
7             if (len(line)!=0):
8                 t=line.split(",")
9                 cases[(t[0],date)]=int(t[1].strip())
10            if (t[0] not in country_list):
11                country_list.append(t[0])
12 death={}
13 for date in date_list:
14     with open("./spider_death/"+date+".csv","r") as f:
15         data=f.readlines()[1:]
16         for line in data:
17             if (len(line)!=0):
18                 t=line.split(",")
19                 death[(t[0],date)]=int(t[1].strip())
20 vaccinated={}
21 for date in date_list:
22     with open("./spider_vaccinated/"+date+".csv","r") as f:
23         data=f.readlines()[1:]
24         for line in data:
25             if (len(line)!=0):
26                 t=line.split(",")
27                 vaccinated[(t[0],date)]=int(t[1].strip())
28 new_cases={}
29 for date in date_list:
30     with open("./spider_new/"+date+".csv","r") as f:
31         data=f.readlines()[1:]
32         for line in data:
33             if (len(line)!=0):
34                 t=line.split(",")
35                 new_cases[(t[0],date)]=int(t[1].strip())
```

4.按要求进行数据分析

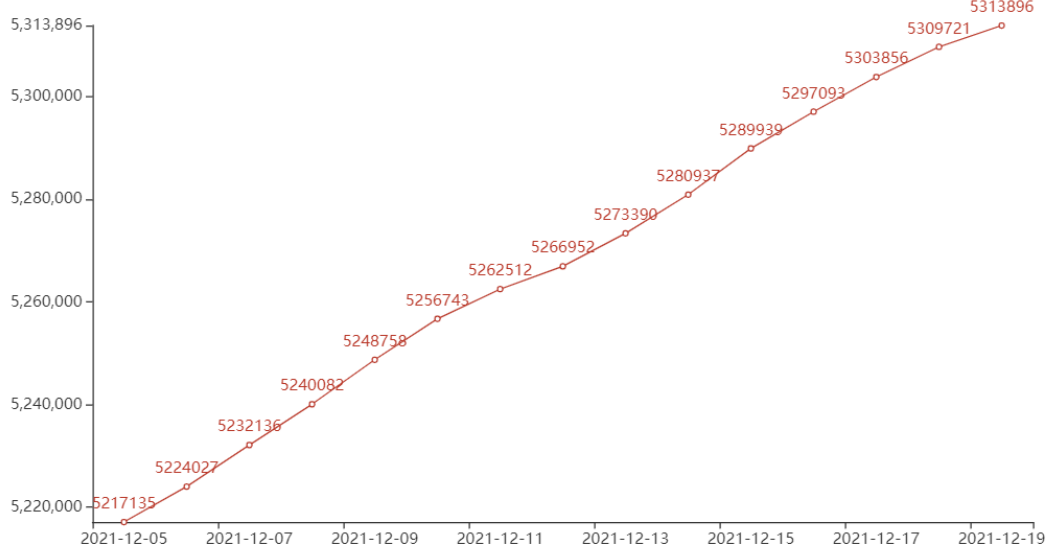
(1)

1) 15 天中，全球新冠疫情的总体变化趋势；

全球确诊总数变化图
(人)



全球死亡病例总数变化图
(人)



- 以上两张图分别从全球总确诊数量，总死亡病例数量出发，分别绘制对应的折线图。
- 可以看出，全球感染数量正在以较为稳定的速度快速增长，疫情仍在快速蔓延。
- 代码：

Python

```

1  # 1
2  #统计确诊总数
3  total_cases=[]
4  for date in date_list:
5      tot=0
6      for country in country_list:
7          tot+=cases[(country,date)]
8      total_cases.append(tot)
9  #导出图表
10 chart = (

```

```

11     Line()
12     .add_xaxis(date_list)
13     .add_yaxis("",total_cases)
14     .set_global_opts(
15         title_opts=opts.TitleOpts(title="全球确诊总数变化图"),
16         yaxis_opts=opts.AxisOpts(
17             type_="value",
18             min_=min(total_cases),
19             max_=max(total_cases)
20         )
21     )
22     .render("全球确诊变化.html")
23 )
24 #统计死亡病例数
25 total_death=[]
26 for date in date_list:
27     tot=0
28     for country in country_list:
29         tot+=death[(country,date)]
30     total_death.append(tot)
31 #导出图表
32 chart = (
33     Line()
34     .add_xaxis(date_list)
35     .add_yaxis("",total_death)
36     .set_global_opts(
37         title_opts=opts.TitleOpts(title="全球死亡病例总数变化图"),
38         yaxis_opts=opts.AxisOpts(
39             type_="value",
40             min_=min(total_death),
41             max_=max(total_death)
42         )
43     )
44     .render("全球死亡病例数变化.html")
45 )

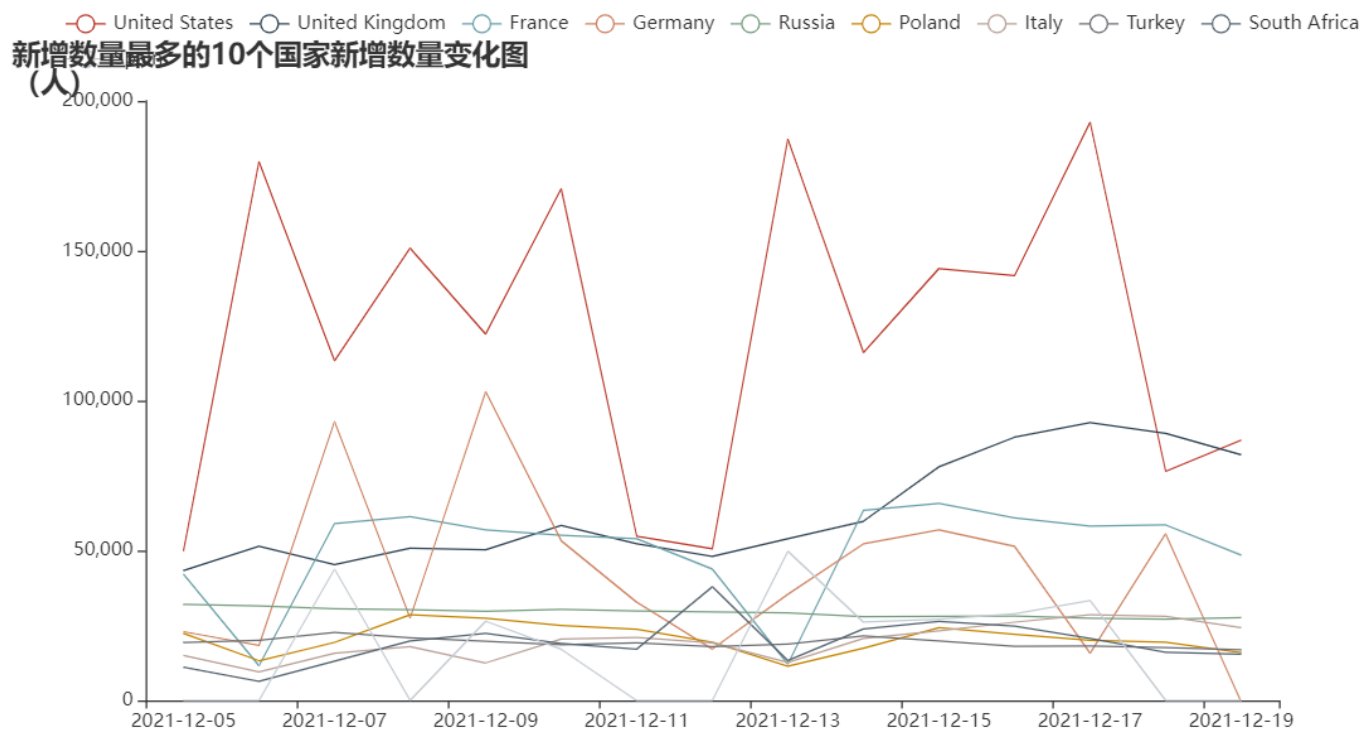
```

(2)

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的曲线图；

首先需要找出这个10个每天新增确诊数量最多的国家，只需要用19日的数据减去5日的数据就好。之后根据每天的新增数量绘制折线图。

- 这新增最多的十个国家如图例所示：



- 可以看出，美国单日新增高居榜首，英法德等东欧国家紧随其后。
- 其次还有南非，俄罗斯，土耳其等区域大国。

代码：

Python

```
1  # 2
2  increase_in_15_days=[]
3  for country in country_list:
4      tot=0
5      for date in date_list:
6          tot+=new_cases[(country,date)]
7      increase_in_15_days.append((tot,country))
8  #对于15天的新增数量进行排序
9  increase_in_15_days.sort(reverse=True)
10 increase_top_10_list=[]
11 #找到新增最多的十个国家
12 for i in range(10):
13     increase_top_10_list.append(increase_in_15_days[i][1])
14 #导出图表
15 chart = (
16     Line()
17     .add_xaxis(date_list)
18     .set_global_opts(
19         title_opts=opts.TitleOpts(title="\n新增数量最多的10个国家新增数量变化图\n\n(人) "),
20         yaxis_opts=opts.AxisOpts(
21             type_="value",
22             min_=0,
23             max_=200000
24         )
25     )
26 )
27 #添加每个国家的折线图
28 for country in increase_top_10_list:
29     increase=[]
30     for date in date_list:
31         t=new_cases[(country,date)]
32         increase.append(t)
33     chart.add_yaxis(country,increase,is_symbol_show=False)
34 chart.render("新增top10.html")
```

(3)

3) 累计确诊数排名前 10 的国家名称及其数量；

- 这个分析比较简单，只需要统计各国确诊数量并排序输出即可。
- 结果如下：

Country	Total Cases
United States	50.89 million
India	34.75 million
Brazil	22.22 million
United Kingdom	11.38 million
Russia	10.04 million
Turkey	9.17 million
France	8.67 million
Germany	6.81 million
Iran	6.17 million
Spain	5.46 million

- 可以发现累计确诊前十名的国家和新增前十名的国家基本差不多
- 还是美国和各个区域的大国。

代码：

Python

```

1  #3
2  total_cases_dict=[]
3  for country in country_list:
4      tot=cases[(country,date_list[-1])]
5      total_cases_dict.append((tot,country))
6  #对确诊数量进行排序
7  total_cases_dict.sort(reverse=True)
8  #导出csv
9  cases_top_10_list=[]
10 with open("./累计确诊top10.csv","w",newline="") as f:
11     writer=csv.writer(f)
12     writer.writerow(["Country","Total Cases"])
13     for i in range(10):
14         t=total_cases_dict[i][0]
15         t/=1000000
16         writer.writerow([total_cases_dict[i][1],str(t)+" million"])

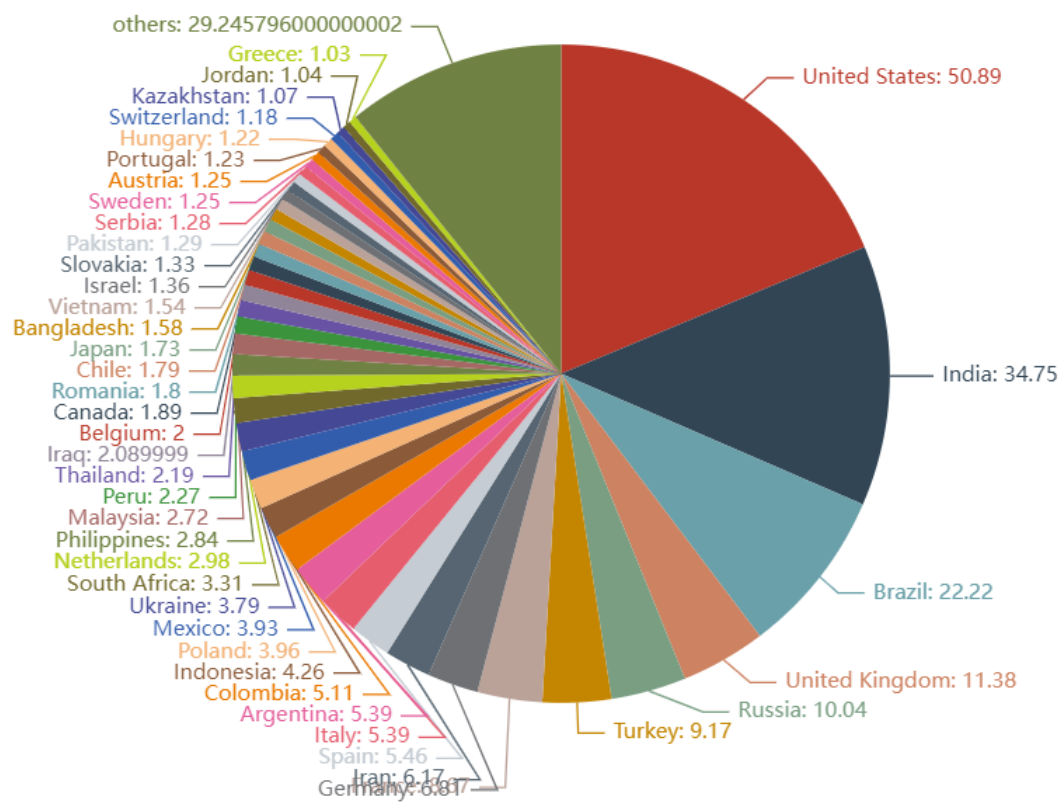
```

(4)

4) 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）；

- 这里对于确诊人数少于一百万的国家合并处理了
- 可以看到美国，印度和巴西总确诊人数占到了全球的三分之一以上，疫情十分严重。

全球累计确诊饼状图 (单位 百万人)



代码:

Python

```
1  # 4
2  #统计各国的疫情确诊比例
3  total_cases_dict_rev=[]
4  others=0
5  for i in total_cases_dict:
6      if (i[0]>1000000):
7          total_cases_dict_rev.append((i[1],i[0]/1000000))
8      else:
9          others+=i[0]/1000000 #小于一百万的合并处理
10 total_cases_dict_rev.append(("others",others))
11 #导出图表
12 chart = (
13     Pie()
14     .add("", total_cases_dict_rev)
15     .set_global_opts(
16         title_opts=opts.TitleOpts(title="全球累计确诊饼状图（单位 百万人）"),
17         legend_opts=opts.LegendOpts(type_="scroll", pos_left="80%", orient="vertical",is_show=False)
18     )
19     .set_series_opts(label_opts=opts.LabelOpts(formatter="{b}: {c}"))
20     .render("全球累计确诊饼状图.html")
21 )
```

(5)

5) 累计确诊人数占国家总人口比例最高的 10 个国家；

- 计算方法是用确诊人数除以国家人口
- 可以发现有的国家确诊比例甚至已经达到四分之一。

Country	Cases Ratio
Andorra	26.564540107297525%
Montenegro	25.633983119259046%
Slovakia	24.355758149885336%
Seychelles	24.312492417195777%
Georgia	22.84456494290492%
Slovenia	21.57871848306942%
San Marino	20.598524267278126%
Lithuania	18.675493389623707%
Estonia	17.516044929575873%
Netherlands	17.352721253164617%

Python

```
1  #5
2  case_ratio=[]
3  case_ratio_country={}
4  for country in country_list:
5      #统计确诊比例
6      tot=cases[(country,date_list[-1])]/population[country]*100
7      case_ratio.append((tot,country))
8      case_ratio_country[country]=tot
9  #对于确诊比例进行排序
10 case_ratio.sort(reverse=True)
11 #导出csv
12 with open("./确诊比例top10.csv","w",newline="") as f:
13     writer=csv.writer(f)
14     writer.writerow(["Country","Cases Ratio"])
15     for i in range(10):
16         t=case_ratio[i][0]
17         writer.writerow([case_ratio[i][1],str(t)+"%"])
```

(6)

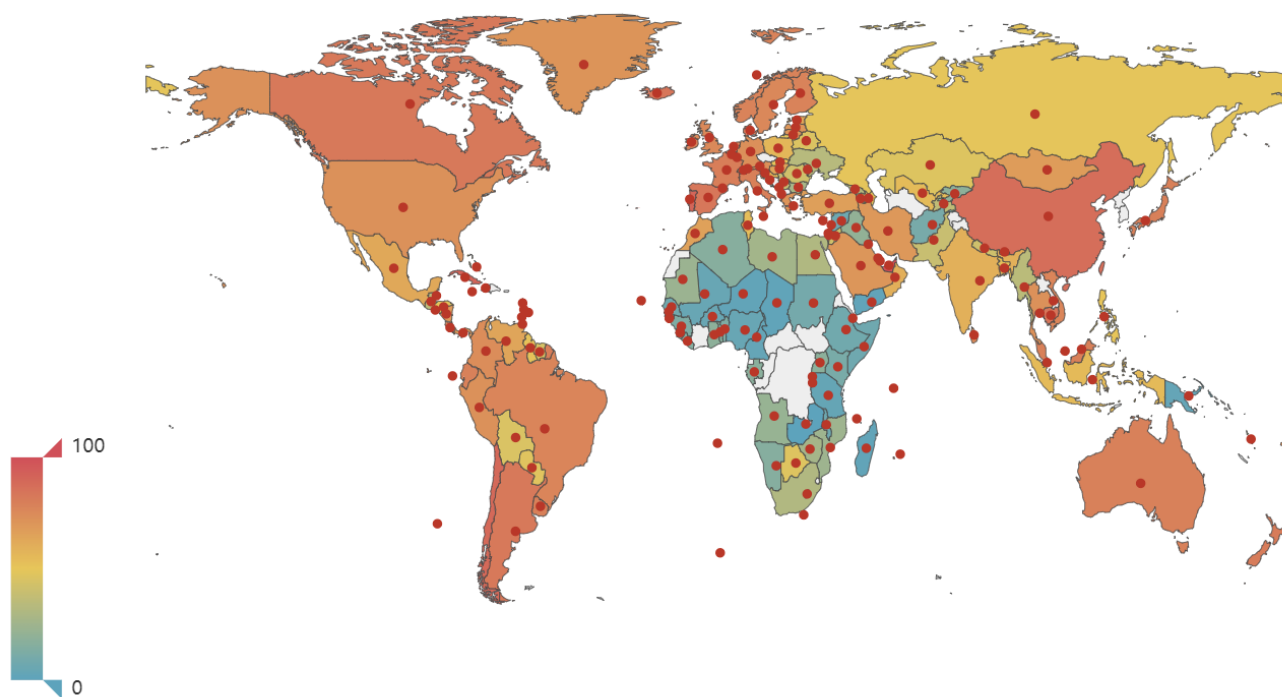
6) 疫苗接种情况（至少接种了一针及以上），请用地图形式展示；

这里使用了pyecharts的地图绘制功能

首先统计各个国家的疫情接种率，按照疫情接种率对各个国家进行染色。

红色表示接种率高，蓝色表示接种率低。

全球疫苗接种率地图



- 可以看出中国接种率很高，美国和欧洲等国家也很好
- 但是非洲国家接种率普遍较低。

Python

```
1  # 6
2  vacc_ratio=[]
3  vacc_map=[]
4  vacc_ratio_country={}
5  #统计每个国家的疫苗接种率
6  for country in country_list:
7      tot=vaccinated[(country,date_list[-1])]/population[country]*100
8      vacc_map.append((country,tot))
9      vacc_ratio.append((tot,country))
10     vacc_ratio_country[country]=tot
11     #导出图表
12     chart = (
13         Map()
14         .add("", vacc_map, "world")
15         .set_series_opts(label_opts=opts.LabelOpts(is_show=False))
16         .set_global_opts(
17             title_opts=opts.TitleOpts(title="全球疫苗接种率地图"),
18             visualmap_opts=opts.VisualMapOpts(max_=100),
19         )
20         .render("全球疫苗接种率地图.html")
21     )
```

(7)

- 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；
- 用统计到的疫苗接种率数据，直接统计最低的10个国家输出即可

Country	Vaccinated Ratio
Burundi	0.03272018214289124%
Haiti	1.06206329491751%
Chad	1.6903473458593075%
Yemen	1.8256487892677884%
South Sudan	1.9967880866446928%
Niger	2.0183307212017816%
Madagascar	2.0719604740902837%
Cameroon	2.96124064322765%
Papua New Guinea	3.1134410423938563%
Tanzania	3.382199778508192%

- 经过验证，这些国家的经济水平较为落后。

Python

```
1 # 7
2 #对接种率进行排序
3 vacc_ratio.sort()
4 #导出csv
5 with open("./接种率last10.csv","w",newline="") as f:
6     writer=csv.writer(f)
7     writer.writerow(["Country","Vaccinated Ratio"])
8     cnt=0
9     i=0
10    while (cnt<10):
11        i+=1
12        t=vacc_ratio[i][0]
13        if (t<1e-5):
14            continue
15        cnt+=1
16        writer.writerow([vacc_ratio[i][1],str(t)+"%"])
```

(8)

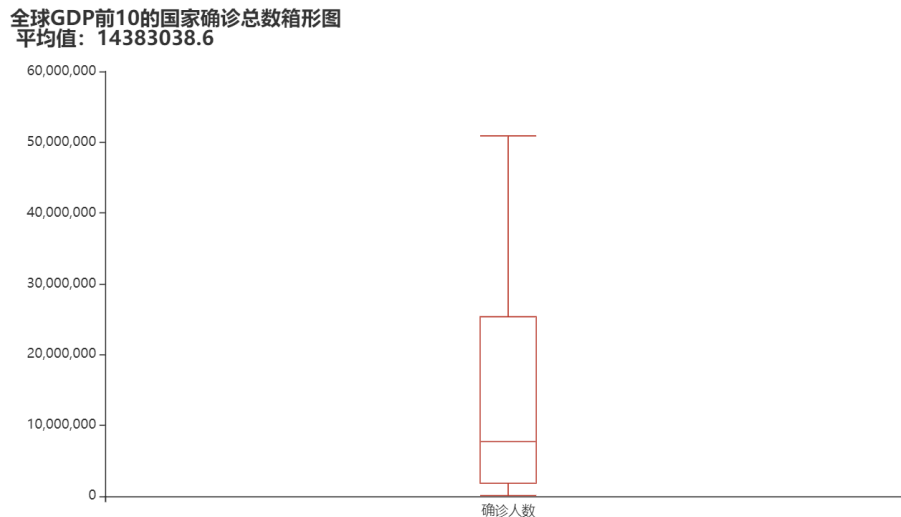
- 8) 全球 GDP 前十名国家的累计确诊人数箱型图，要有平均值；

这需要先查出GDP前十名的国家有哪些：

Python

```
1 gdp_top_10_country=["United States","China","Japan","Germany","India","United Kingdom","France","Brazil","Italy","Canada"]
```

之后统计这些国家的确诊人数，用pyecharts绘制箱型图，同时计算平均值并输出。



- 可以看出整体数据跨度非常大，最低的是中国，最高的是美国。
- 平均值是一千四百万人，不过中位数不到一千万，平均值是由美国拉上去的。
- 代码：

Python

```
1 # 8
2 gdp_top_10_country=["United States","China","Japan","Germany","India","United Kingdom","France","Brazil","Italy","Canada"]
3 box_data=[]
4 #统计这些国家的确诊总数
5 for country in gdp_top_10_country:
6     box_data.append(cases[(country,date_list[-1])])
7 #导出图表
8 c = Boxplot()
9 c.add_xaxis(["确诊人数"])
10 c.add_yaxis("", c.prepare_data([box_data]))
11 c.set_global_opts(title_opts=opts.TitleOpts(title="全球GDP前10的国家确诊总数箱形图\n 平均值: {}".format(sum(box_data)/len(box_data))))
12 c.render("GDP前10的国家确诊总数箱形图.html")
```

(9)

9) 死亡率最高的 10 个国家；

用死亡病例数字除以总人数，就得到死亡率，由此统计即可。

Country	Death Ratio
Peru	26.564540107297525%
Bulgaria	25.633983119259046%
Bosnia and Herz	24.355758149885336%
Hungary	24.312492417195777%
Montenegro	22.84456494290492%
North Macedonia	21.57871848306942%
Georgia	20.598524267278126%
Romania	18.675493389623707%
Croatia	17.516044929575873%
Slovakia	17.352721253164617%

代码：

Python

```

1  # 9
2  death_ratio=[]
3  death_ratio_country={}
4  for country in country_list:
5      tot=death[(country,date_list[-1])]/population[country]*100
6      death_ratio.append((tot,country))
7      death_ratio_country[country]=tot
8  #按照死亡率进行排序
9  death_ratio.sort(reverse=True)
10 #导出csv
11 with open("./死亡比例top10.csv","w",newline="") as f:
12     writer=csv.writer(f)
13     writer.writerow(["Country","Death Ratio"])
14     for i in range(10):
15         t=case_ratio[i][0]
16         writer.writerow([death_ratio[i][1],str(t)+"%"])

```

可以看出死亡率较高的国家仍然是一些疫苗接种率低，并且经济发展水平低的国家。

(10)

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。

考虑已有数据，可以通过综合死亡率，确诊率，疫苗接种情况来综合设计一个评分，来评估一个国家的防疫现状。

我认为一个国家的死亡率越低，确诊率越低，疫苗接种率越高，可以说明这个国家应对疫情最好，所以对死亡率，确诊率和疫苗接种率加权处理得到评分。

我设计的一个评分标准是： $score = \text{疫苗接种率} - 5 * \text{确诊比例} - 5 * \text{死亡率} + 50$

这个标准不一定很科学，仅作为尝试参考。

这样得到一个分数，取分数最高的十个国家：

Country	Score
China	274.42
Cambodia	260.66
New Zealand	260.15
South Korea	258.21
Taiwan	255.2
Brunei	248.39
Bhutan	247.83
Australia	247.71
Japan	245.11
Vietnam	238.24

代码：

Python

```
1 score=[]
2 score_map=[]
3 #对每个国家计算分数
4 for country in country_list:
5     s=50+vacc_ratio_country[country]-5*case_ratio_country[country]-5*death_ratio_country[country]
6     score.append((s, country))
7     score_map.append((country, 275-s))
8 score.sort(reverse=True)
9 #导出csv
10 with open("./抗疫最好的top10.csv", "w", newline="") as f:
11     writer=csv.writer(f)
12     writer.writerow(["Country", "Score"])
13     for i in range(10):
14         t=score[i][0]
15         writer.writerow([score[i][1], str(t)])
```

(11)

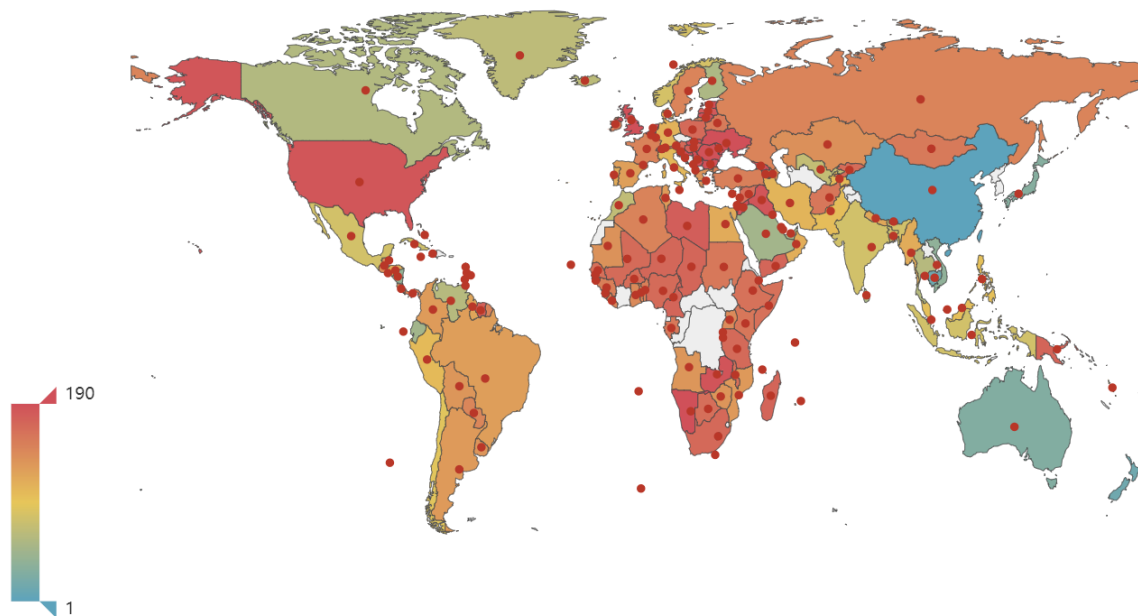
10) 其它你希望分析和展示的数据。

我选择了两个角度：

- 根据上一题设计的评分标准，分析并绘制全球国家抗疫水平地图：
- 绘制确诊比例和疫苗接种率的散点图，分析其关联

全球国家抗疫水平地图

全球抗疫地图



从这张图中可以显著看出，中国防疫取得了非常大的成果
而美国和欧洲防疫状况堪忧。

代码：

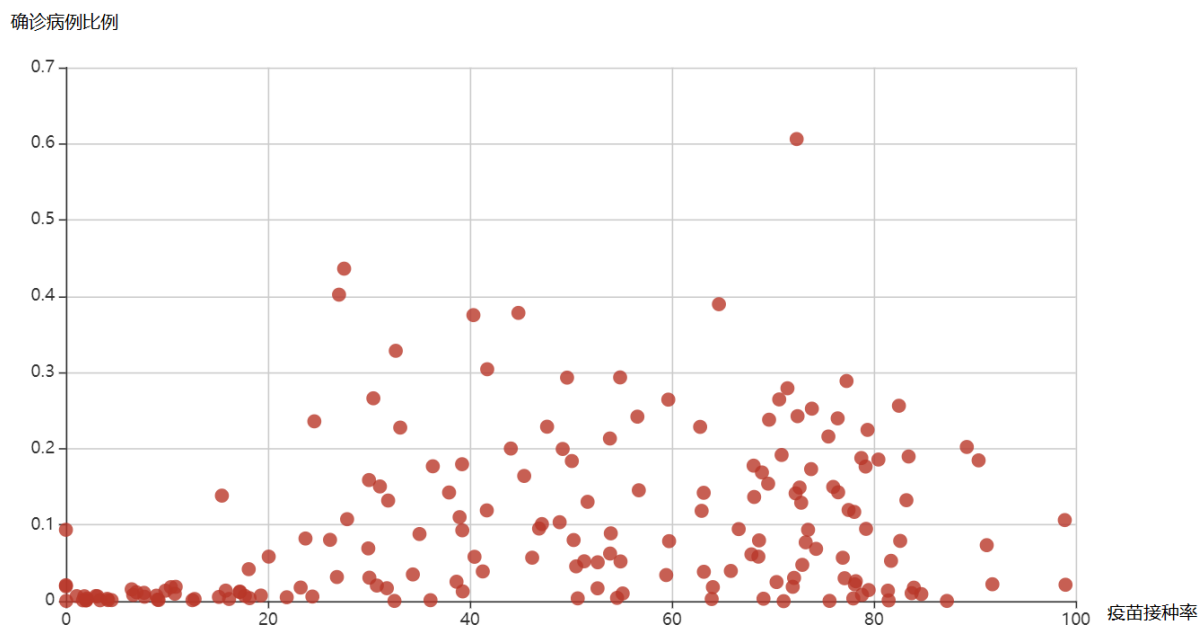
Python

```
1 for i in score_map:
2     maxx=max(maxx,i[1])
3     minn=min(minn,i[1])
4 #导出图表
5 chart = (
6     Map()
7     .add("", score_map, "world")
8     .set_series_opts(label_opts=opts.LabelOpts(is_show=False))
9     .set_global_opts(
10         title_opts=opts.TitleOpts(title="全球抗疫地图"),
11         visualmap_opts=opts.VisualMapOpts(max_=190,min_=minn),
12     )
13     .render("全球抗疫地图.html")
14 )
```

确诊比例和疫苗接种率散点图

为了研究确诊比例和疫苗接种率之间的关系，我将这所有的国家的确诊比例和疫苗接种率绘制成散点图，纵轴是确诊比例。横轴是疫苗接种率。

得到结果如下：



可以发现，疫苗接种率和确诊病例的比例并没有十分明确的关系。

分析：

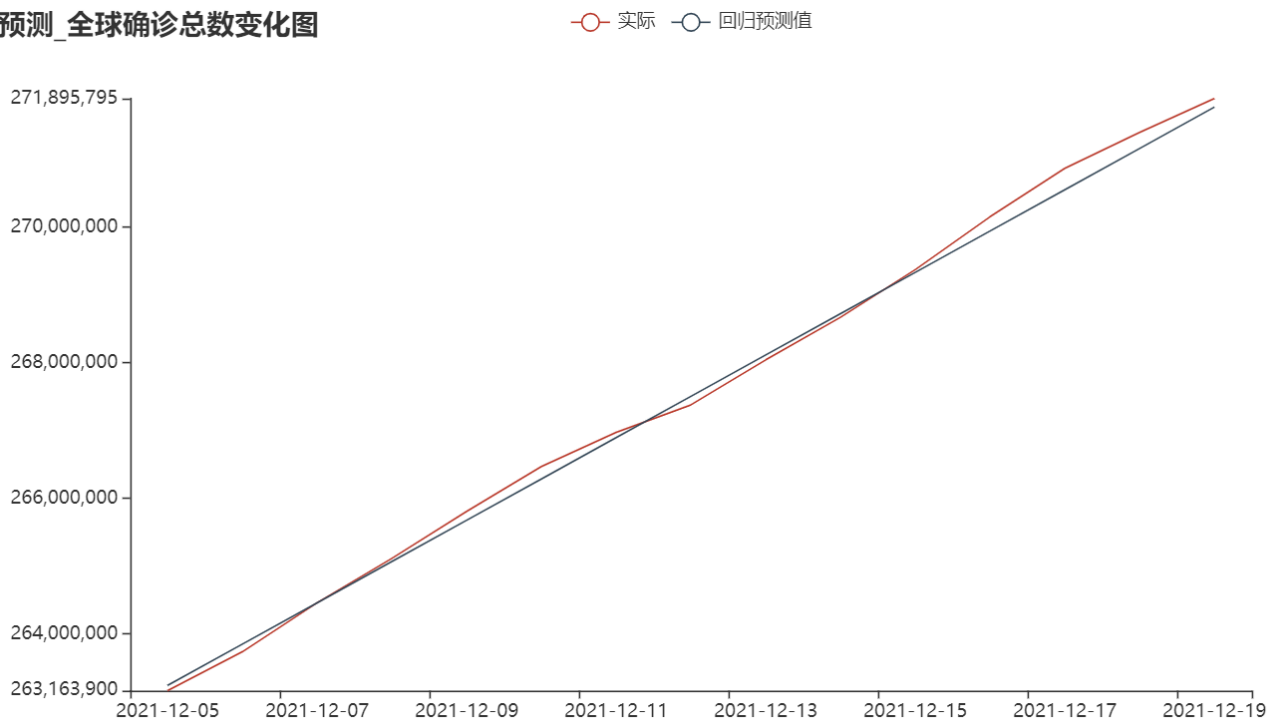
- 疫苗接种可以部分算作是疫情蔓延的结果，可能很多国家是因为疫情严重，其他防控措施无效，大力推广疫苗，所以并没有预期的疫苗接种率越高，患病率越低的关系。
- 相比起疫苗，政府的防控措施，民众意识，地理位置，医疗水平，气候等因素对于疫情发展的影响也非常大，所以疫苗接种的影响不容易直观的看出来
- 新冠病毒有很多如奥密克戎的变种，已经接种的疫苗对于各种变种的免疫力有所下降甚至消失，这也是疫苗作用不明显的原因。

(12)

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。

- 观察全球确诊病例总数变化图，发现增长很稳定，接近线性函数
- 于是尝试使用线性回归进行预测
- 使用了scipy库的线性回归函数
- 结果如下：

预测_全球确诊总数变化图



其中，红线表示的是实际的确诊人数，而灰色的线表示的是回归的曲线。

这条灰色的曲线是由前十天的数据计算得到的，可以看出对于后五天的数据预测结果还不错。

但是仍需意识到，疫情的传播速度和自然状况，基因突变，各国政策息息相关，所以很难用一个简单的方法来预测疫情的发展。

四、实验总结

经过这次实验，我学会了使用python进行简单的数据分析，对于python常用库和函数的使用也更加熟练，对于Python的解释器也有了更加深入的理解。这次实验我从多个方面入手，通过爬虫爬取的数据，加以整理。最终以图表的形式呈现出来。我深深感受到了数据分析和统计的重要性，也感受到了python这个工具对于数据分析所带来的便利性和便捷性。

这次统计疫情数据，让我意识到全球抗疫之路前路漫漫，让我意识到相对于疫情肆虐的西方国家，中国的抗疫举措多么有效，抗疫成果多么显著。我深深的为自己是一个中国人而感到自豪。

在今后的学习生活中，我将注重自己对于信息获取数据分析方面技能的培养，将这种数据分析的思维和意识融入到生活中。同时进一步增强对于各种工具的使用的学习，学会从各个角度全面的思考和分析问题。