

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sn
5
```

背景介绍
新型冠状病毒感染的肺炎疫情爆发后，对人们的生活产生很大的影响。
当前感染人数依然在不断变化。每天国家卫健委和各大新闻媒体都会公布疫情的数据，包括累计确诊人数、现有确诊人数等。
本案例使用python对新冠肺炎的数据（韩国）进行分析。

读取数据集

```
1 patientdata=pd.read_csv("D:360Downloads\PatientInfo.csv")
2 region=pd.read_csv("D:360Downloads\Region.csv")
```

Task1

数据预处理

查看数据前五

```
1 patientdata.head()
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	patient_id	global_num	sex	birth_year	age	country	province	city	disease	infection_case	info
0	1000000001	2.0	male	1964.0	50s	Korea	Seoul	Gangseo-gu	NaN	overseas inflow	1.0
1	1000000002	5.0	male	1987.0	30s	Korea	Seoul	Jungnang-gu	NaN	overseas inflow	1.0
2	1000000003	6.0	male	1964.0	50s	Korea	Seoul	Jongno-gu	NaN	contact with patient	2.0
3	1000000004	7.0	male	1991.0	20s	Korea	Seoul	Mapo-gu	NaN	overseas inflow	1.0
4	1000000005	9.0	female	1992.0	20s	Korea	Seoul	Seongbuk-gu	NaN	contact with patient	2.0

检测数据缺失值

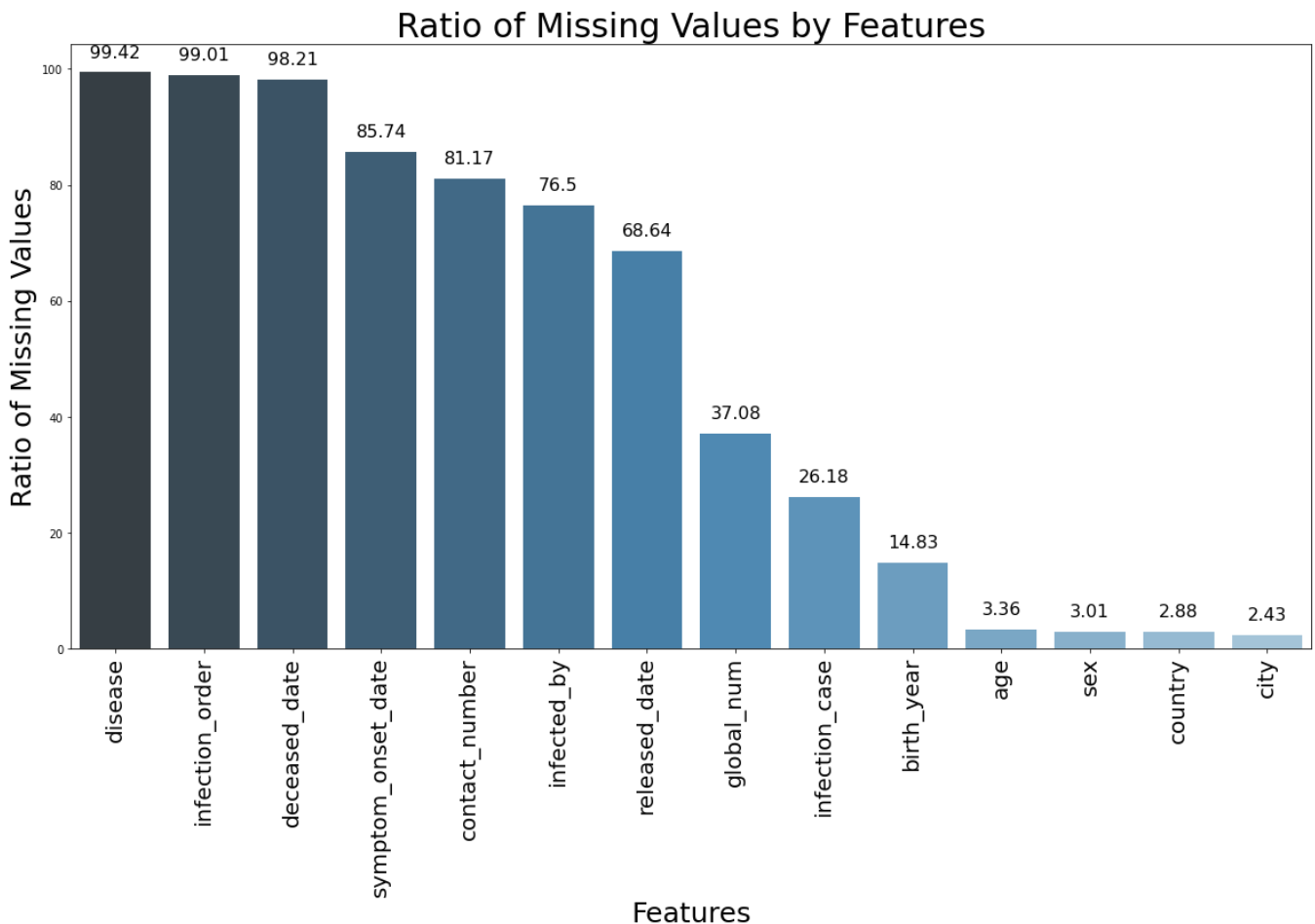
```
1 patientdata.isnull().sum()
2
```

```
1 patient_id      0
2 global_num      1160
3 sex             94
4 birth_year      464
5 age             105
6 country         90
7 province        0
8 city           76
9 disease         3110
10 infection_case  819
11 infection_order 3097
12 infected_by     2393
13 contact_number  2539
14 symptom_onset_date 2682
15 confirmed_date  0
16 released_date   2147
17 deceased_date   3072
18 state          0
```

```
19 | dtype: int64
```

```
1 | na_ratio = pd.DataFrame(patientdata.isnull().sum()/len(patientdata)*100,columns=['NA_Ratio'])
2 | p_na = na_ratio[na_ratio['NA_Ratio']>0].sort_values(by = 'NA_Ratio', ascending=False)
```

```
1 | fig, ax = plt.subplots(figsize=(20,10))
2 | nar = sn.barplot(x=p_na.index, y=p_na['NA_Ratio'],orient='v',palette="Blues_d")
3 | ax.set_xticklabels(p_na.index, rotation=90, fontsize=20)
4 | ax.set_ylabel('Ratio of Missing Values', fontsize=25)
5 | ax.set_xlabel('Features', fontsize=25)
6 | ax.set_title('Ratio of Missing Values by Features', fontsize=30)
7 | for loc, value in zip(ax.patches, p_na.NA_Ratio):
8 |     ax.text(loc.get_x()+loc.get_width()/2, loc.get_height()+2, round(value,2),ha='center', va='bottom',fontsize=16)
9 | plt.show()
```



我们对于图表当中缺失的特征值进行处理

我们对于'disease','infected_by','symptom_onset_date','deceased_date','released_date',NaN表示没有这个特征，可以将NaN用None代替。

```
1 | none_col=['disease', 'infected_by','symptom_onset_date','deceased_date','released_date']
2 | patientdata[none_col]=patientdata[none_col].fillna('None')
```

对于'contact_number'用平均值填补

```
1 | patientdata['contact_number'].fillna(patientdata['contact_number'].mean(),inplace=True)
```

对于'country','city','province','infection_case','age','sex','birth_year'采用众数填补

```
1 | patientdata['city'].fillna(patientdata['city'].mode()[0],inplace=True)
2 | patientdata['country'].fillna(patientdata['country'].mode()[0],inplace=True)
3 | patientdata['infection_case'].fillna(patientdata['infection_case'].mode()[0],inplace=True)
4 | patientdata['age'].fillna(patientdata['age'].mode()[0],inplace=True)
5 | patientdata['sex'].fillna(patientdata['sex'].mode()[0],inplace=True)
6 | patientdata['birth_year'].fillna(patientdata['birth_year'].mode()[0],inplace=True)
7 |
```

去除无用特征ID,因为global_num以及infection_order缺失的比例过大，考虑删除这两个特征,因为age和birthyear有多重共线性，只考虑birthyear，删除age

```

1 patientdata.drop('patient_id',axis=1,inplace=True)
2 patientdata.drop('global_num',axis=1,inplace=True)
3 patientdata.drop('infection_order',axis=1,inplace=True)
4 patientdata.drop('age',axis=1,inplace=True)
5

```

再次查看缺失，确认已完成缺失值处理

```

1 patientdata.isnull().sum()

```

```

1 sex                0
2 birth_year        0
3 country           0
4 province          0
5 city              0
6 disease           0
7 infection_case     0
8 infected_by       0
9 contact_number    0
10 symptom_onset_date 0
11 confirmed_date    0
12 released_date     0
13 deceased_date     0
14 state            0
15 dtype: int64

```

region中，将kindergarden， elementary_school_count,university_count加在一起，用education_healthy_facilities表示

```

1 region['education_healthy_facilities']=region['kindergarten_count']+region['elementary_school_count']+region['university_count']
2 region.head()

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	code	province	city	latitude	longitude	elementary_school_count	kindergarten_count	university_count	academy_ratio
0	10000	Seoul	Seoul	37.566953	126.977977	607	830	48	1.44
1	10010	Seoul	Gangnam-gu	37.518421	127.047222	33	38	0	4.18
2	10020	Seoul	Gangdong-gu	37.530492	127.123837	27	32	0	1.54
3	10030	Seoul	Gangbuk-gu	37.639938	127.025508	14	21	0	0.67
4	10040	Seoul	Gangseo-gu	37.551166	126.849506	36	56	1	1.17

将region的数据导入patientdata中，作为patient新的特征列

```

1
2 def getRegion(dataframe,dR):
3
4     dataframe['education_healthy_facilities']=0
5     dataframe['latitude']=0
6     dataframe['longitude']=0
7     dataframe['academy_ratio']=0
8     dataframe['elderly_population_ratio']=0
9     for i in dataframe.index:
10         for j in dR.index:
11
12             if dataframe.loc[i,'city']==dR.loc[j,'city']:
13                 dataframe.loc[i,'latitude']=dR.loc[j,'latitude']
14                 dataframe.loc[i,'longitude']=dR.loc[j,'longitude']
15                 dataframe.loc[i,'education_healthy_facilities']=dR.loc[j,'education_healthy_facilities']
16                 dataframe.loc[i,'academy_ratio']=dR.loc[j,'academy_ratio']
17                 dataframe.loc[i,'elderly_population_ratio']=dR.loc[j,'elderly_population_ratio']
18

```

```
19 | return dataframe
```

```
1 | getRegion(patientdata,region)
2 |
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

	sex	birth_year	country	province	city	disease	infection_case	infected_by	contact_number	symptom_onset_da
0	male	1964.0	Korea	Seoul	Gangseo-gu	None	overseas inflow	None	75.0	2020-01-22
1	male	1987.0	Korea	Seoul	Jungnang-gu	None	overseas inflow	None	31.0	None
2	male	1964.0	Korea	Seoul	Jongno-gu	None	contact with patient	2.002e+09	17.0	None
3	male	1991.0	Korea	Seoul	Mapo-gu	None	overseas inflow	None	9.0	2020-01-26
4	female	1992.0	Korea	Seoul	Seongbuk-gu	None	contact with patient	1e+09	2.0	None
...
3123	female	1995.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	20.0	None
3124	male	1995.0	United States	Jeju-do	Jeju-do	None	overseas inflow	None	23.0	None
3125	female	1996.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	26.0	None
3126	female	1995.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	25.0	None
3127	female	1995.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	14.0	None

3128 rows × 19 columns

查看特征类型

```
1 | type_p=patientdata.columns.groupby(patientdata.dtypes)
2 | [(i,len(type_p[i]),type_p[i]) for i in type_p]
```

```
1 | [(dtype('int64'), 1, Index(['education_healthy_facilities'], dtype='object')),
2 |  (dtype('float64'),
3 |  6,
4 |  Index(['birth_year', 'contact_number', 'latitude', 'longitude',
5 |        'academy_ratio', 'elderly_population_ratio'],
6 |        dtype='object')),
7 |  (dtype('o'),
8 |  12,
9 |  Index(['sex', 'country', 'province', 'city', 'disease', 'infection_case',
10 |        'infected_by', 'symptom_onset_date', 'confirmed_date', 'released_date',
11 |        'deceased_date', 'state'],
12 |        dtype='object'))]
```

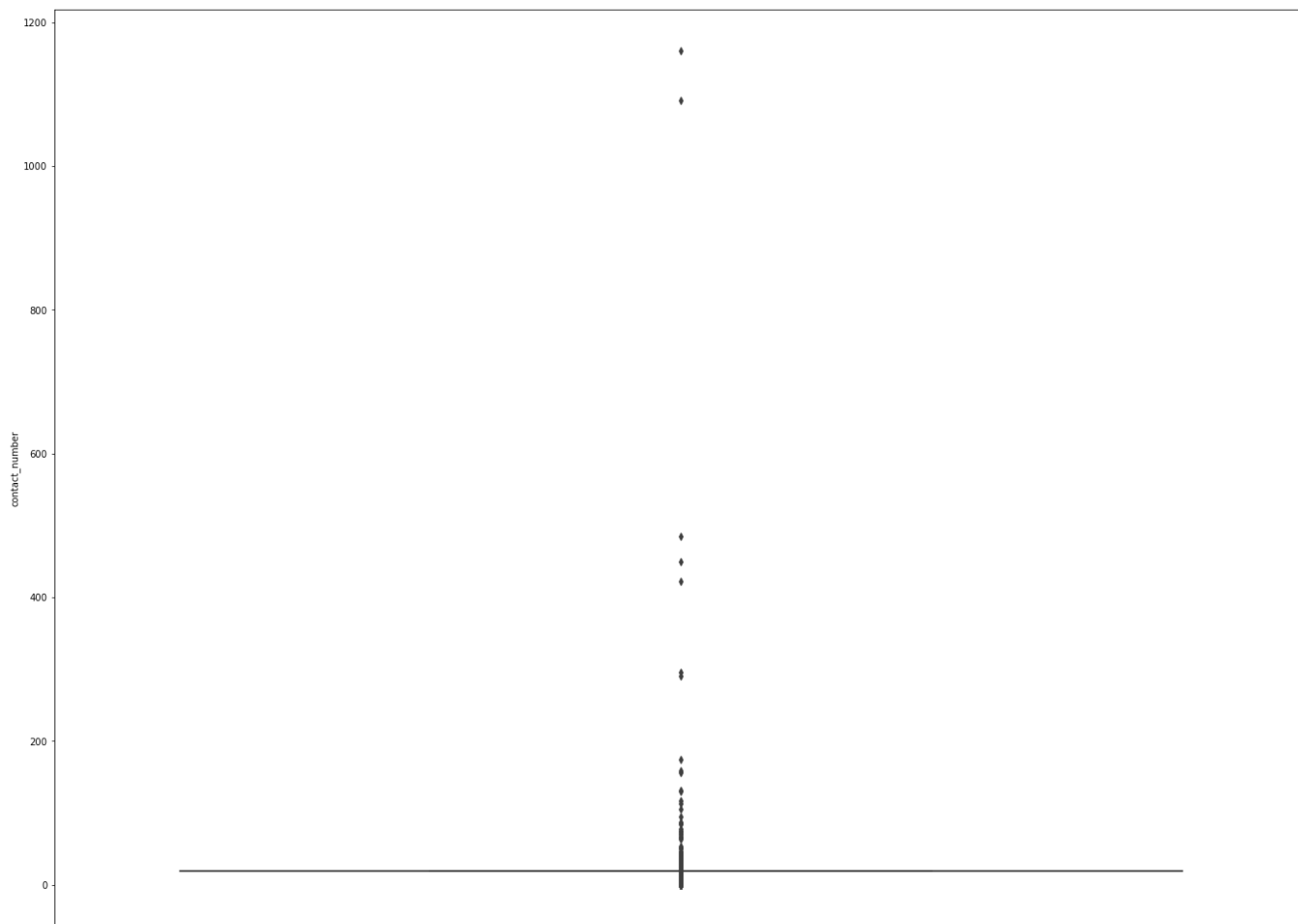
可以知道数值特征有6个，非数值特征有12个

首先对数值型特征去除异常值

首先，对异常值进行查看

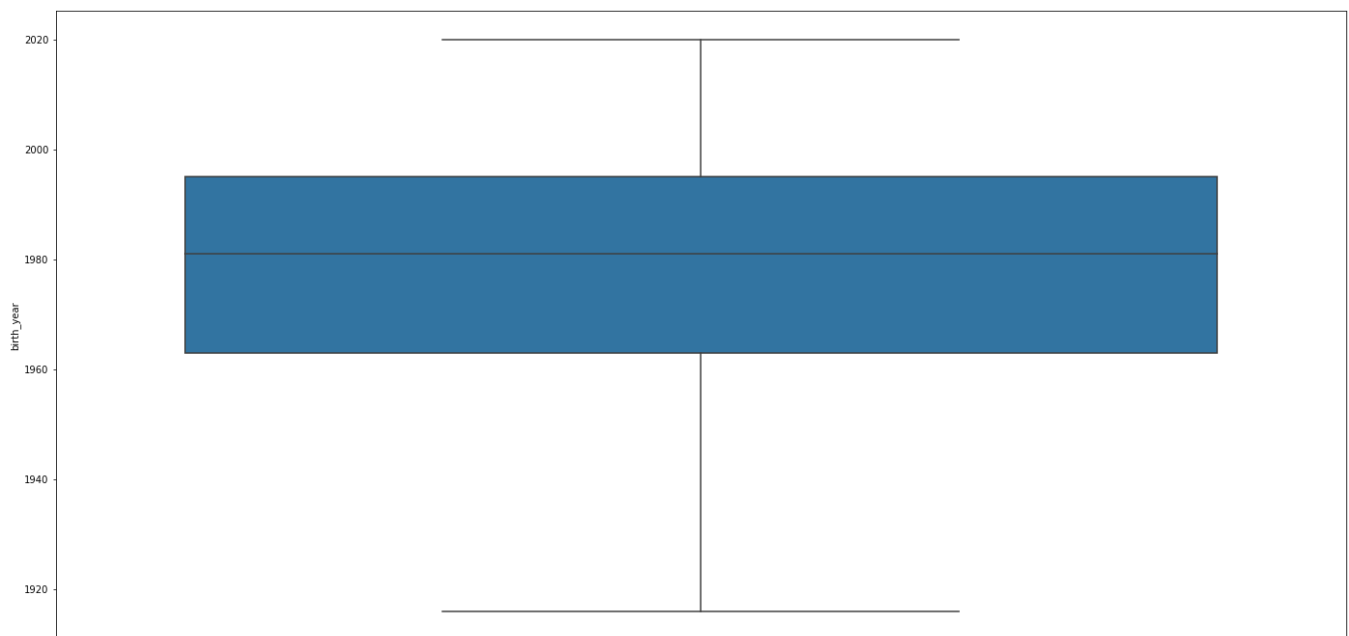
```
1 | fig, axes = plt.subplots(1,1)
2 | fig.set_size_inches(24,18)
3 | sn.boxplot(data=patientdata, y="contact_number")
4 |
```

```
1 | <matplotlib.axes._subplots.AxesSubplot at 0x180cb7ea2e0>
```



```
1 fig, axes = plt.subplots(1,1)
2 fig.set_size_inches(24,12)
3 sn.boxplot(data=patientdata, y="birth_year")
4
```

```
1 | <matplotlib.axes._subplots.AxesSubplot at 0x180cd9ac490>
```



删除小于contact_number小于0, 以及birthyear大于2020, 小于1920的数据

```
1 patientdata.drop(patientdata[patientdata['contact_number']<0].index)
2 patientdata.drop(patientdata[(patientdata['birth_year']<1920)|(patientdata['birth_year']>2020)].index)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	sex	birth_year	country	province	city	disease	infection_case	infected_by	contact_number	symptom_onset_da
0	male	1964.0	Korea	Seoul	Gangseo-gu	None	overseas inflow	None	75.0	2020-01-22
1	male	1987.0	Korea	Seoul	Jungnang-gu	None	overseas inflow	None	31.0	None
2	male	1964.0	Korea	Seoul	Jongno-gu	None	contact with patient	2.002e+09	17.0	None
3	male	1991.0	Korea	Seoul	Mapo-gu	None	overseas inflow	None	9.0	2020-01-26
4	female	1992.0	Korea	Seoul	Seongbuk-gu	None	contact with patient	1e+09	2.0	None
...
3123	female	1995.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	20.0	None
3124	male	1995.0	United States	Jeju-do	Jeju-do	None	overseas inflow	None	23.0	None
3125	female	1996.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	26.0	None
3126	female	1995.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	25.0	None
3127	female	1995.0	Korea	Jeju-do	Jeju-do	None	overseas inflow	None	14.0	None

3127 rows × 19 columns

数字编码, 使用sklearn中的LabelEncoder对有优劣之分或次序关系的分类特征进行数字编码

```
1 #导入LabelEncoder
2 from sklearn.preprocessing import LabelEncoder
```

```

1 column=['sex', 'country', 'province', 'city', 'infection_case', 'state']
2 for col in column:
3     le=LabelEncoder()
4     patientdata[col]=le.fit_transform(patientdata[col])
5

```

对patientinfo的各个特征进行可视化处理

```

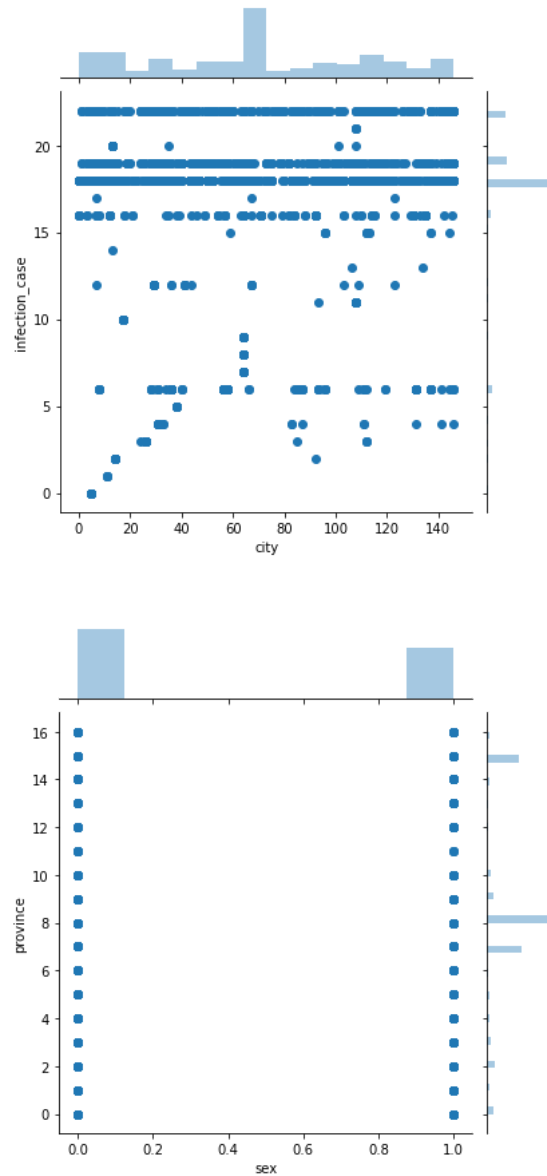
1 sn.jointplot(x="city", y="infection_case", data=patientdata, kind="scatter")
2 sn.jointplot(x="sex", y="province", data=patientdata, kind="scatter")
3 sn.jointplot(x="infection_case", y="state", data=patientdata, kind="scatter")
4 sn.jointplot(x="country", y="state", data=patientdata, kind="scatter")
5 sn.jointplot(x="birth_year", y="state", data=patientdata, kind="scatter")

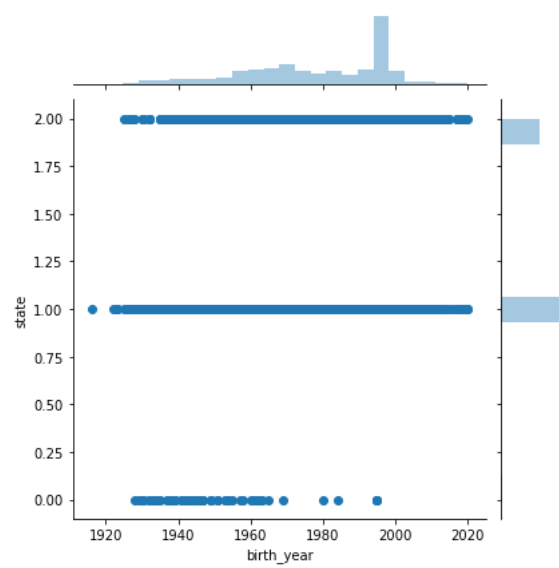
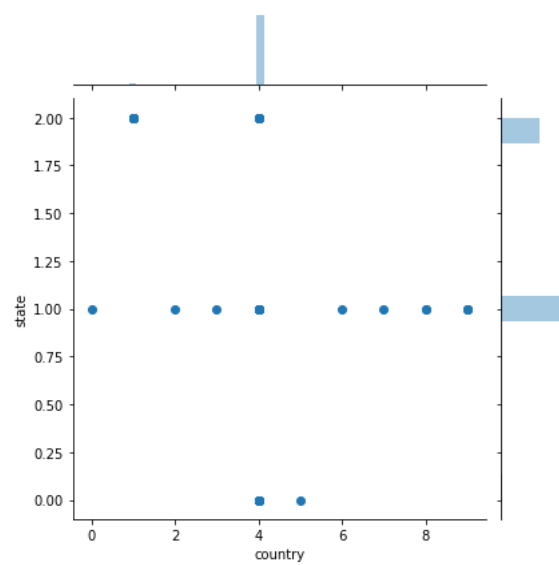
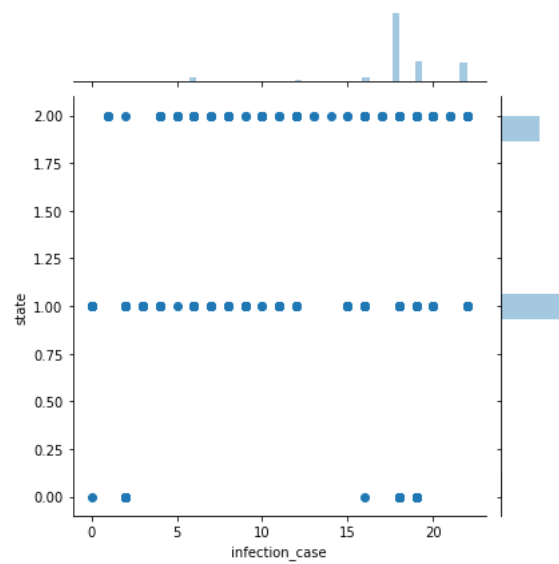
```

```

1 <seaborn.axisgrid.JointGrid at 0x180cdd47100>

```





1 | # 对 'time', 'Timegender', 'Timeage', 'Timeprovince' 进行可视化处理

读取数据集'Time','Timegender','Timeage','Timeprovince'


```

1 time=pd.read_csv("D:360Downloads\Time.csv")
2 timeage=pd.read_csv("D:360Downloads\TimeAge.csv")
3 timegender=pd.read_csv("D:360Downloads\TimeGender.csv")
4 timeprovince=pd.read_csv("D:360Downloads\TimeProvince.csv")

```

```

1 time.head()

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

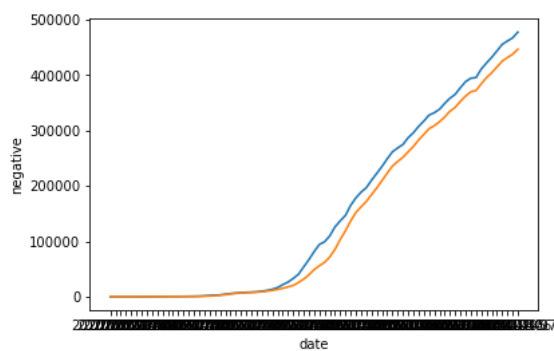
	date	time	test	negative	confirmed	released	deceased
0	2020-01-20	16	1	0	1	0	0
1	2020-01-21	16	1	0	1	0	0
2	2020-01-22	16	4	3	1	0	0
3	2020-01-23	16	22	21	1	0	0
4	2020-01-24	16	27	25	2	0	0

对于time文件里的各项属性进行可视化处理

```

1 ax=sn.lineplot(x='date',y='test',data=time)
2 ax=sn.lineplot(x='date',y='negative',data=time)
3

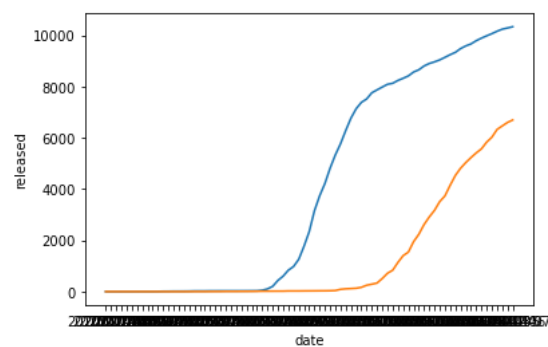
```



```

1 ax=sn.lineplot(x='date',y='confirmed',data=time)
2 ax=sn.lineplot(x='date',y='released',data=time)
3

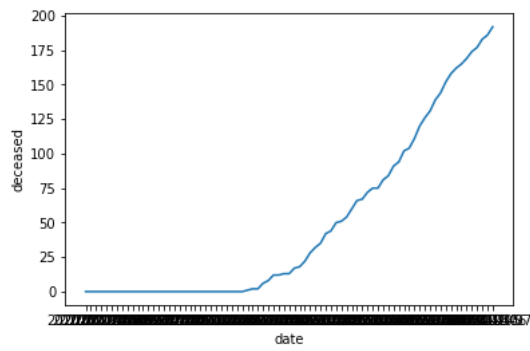
```



```

1 ax=sn.lineplot(x='date',y='deceased',data=time)

```



再对timeage的数据进行可视化处理

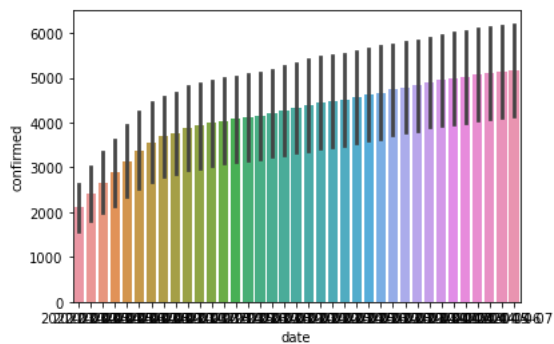
```
1 | timeage.head()
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

	date	time	age	confirmed	deceased
0	2020-03-02	0	0s	32	0
1	2020-03-02	0	10s	169	0
2	2020-03-02	0	20s	1235	0
3	2020-03-02	0	30s	506	1
4	2020-03-02	0	40s	633	1

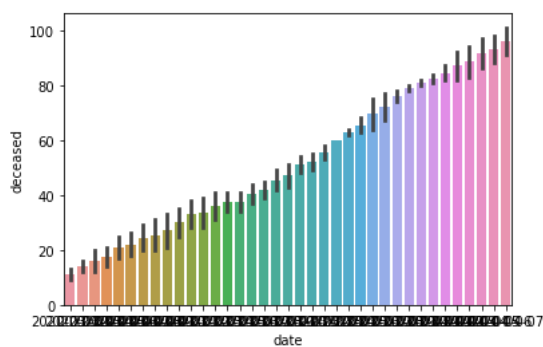
对'confirmed'可视化

```
1 | # age = timeage['age']
2 | # ax =sn.distplot(age)
3 | ax=sn.barplot(x='date',y='confirmed',data=timegender)
```



对deceased可视化

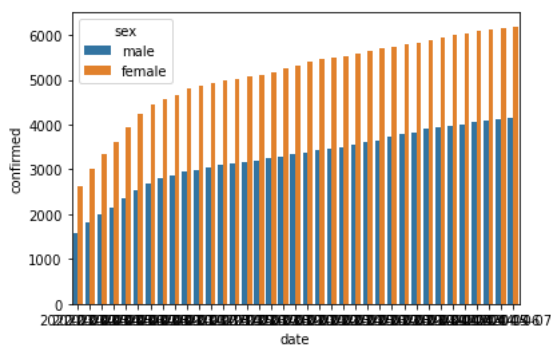
```
1 | # deceased = timeage['deceased']
2 | # ax =sn.distplot(deceased)
3 | ax=sn.barplot(x='date',y='deceased',data=timegender)
```



对timegender进行可视化处理

对'confirmed'和'sex'进行关联可视化

```
1 # time = timegender['sex']
2 # ax =sn.distplot(time)
3 ax=sn.barplot(x='date',y='confirmed',hue = 'sex',data=timegender)
```

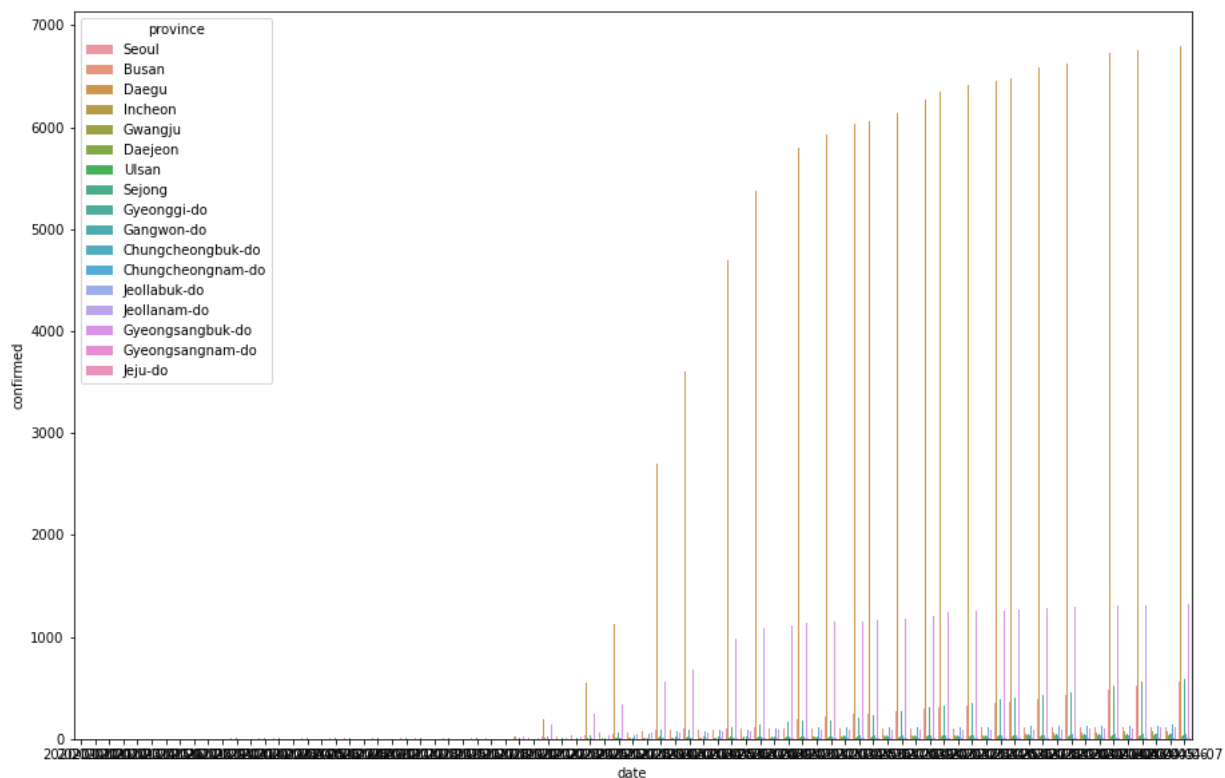


从上表中看出女性确诊数要多于男性

对timeprovince进行可视化处理

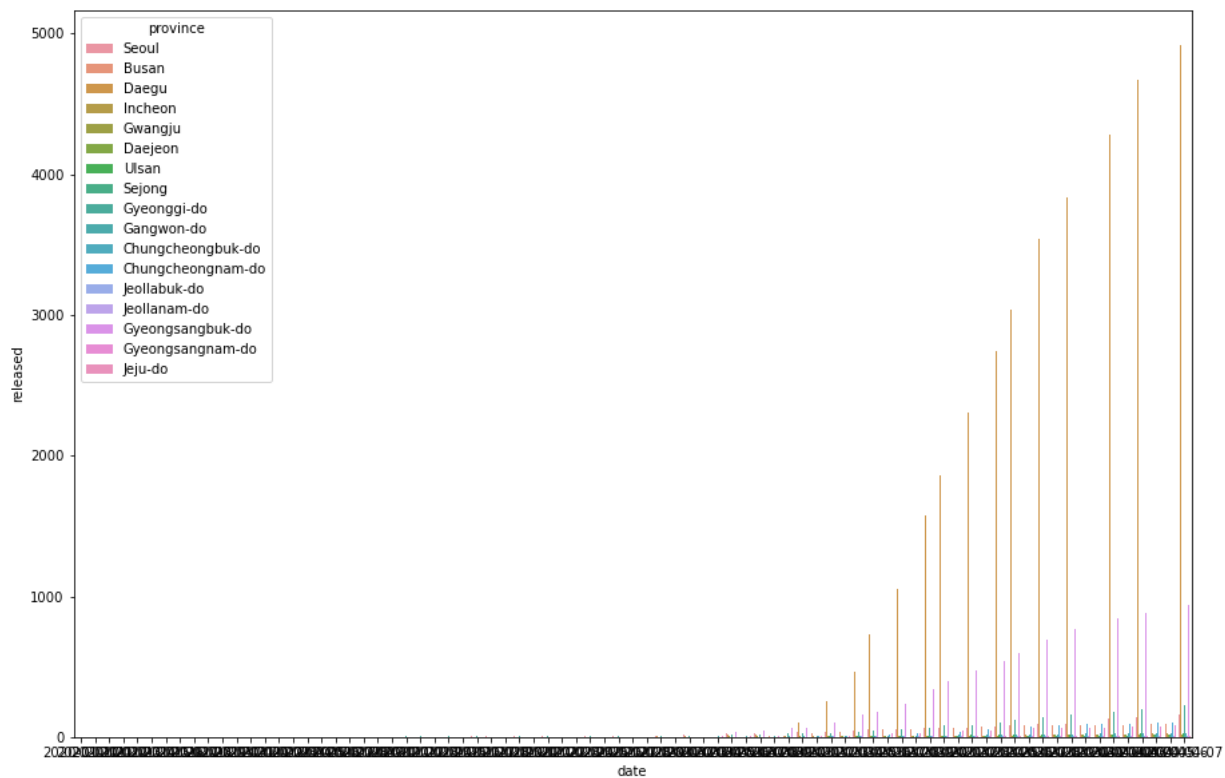
对confirmed和province进行关联可视化处理

```
1 plt.figure(figsize=(15, 10))
2 ax=sn.barplot(x='date',y='confirmed',hue = 'province',data=timeprovince)
```



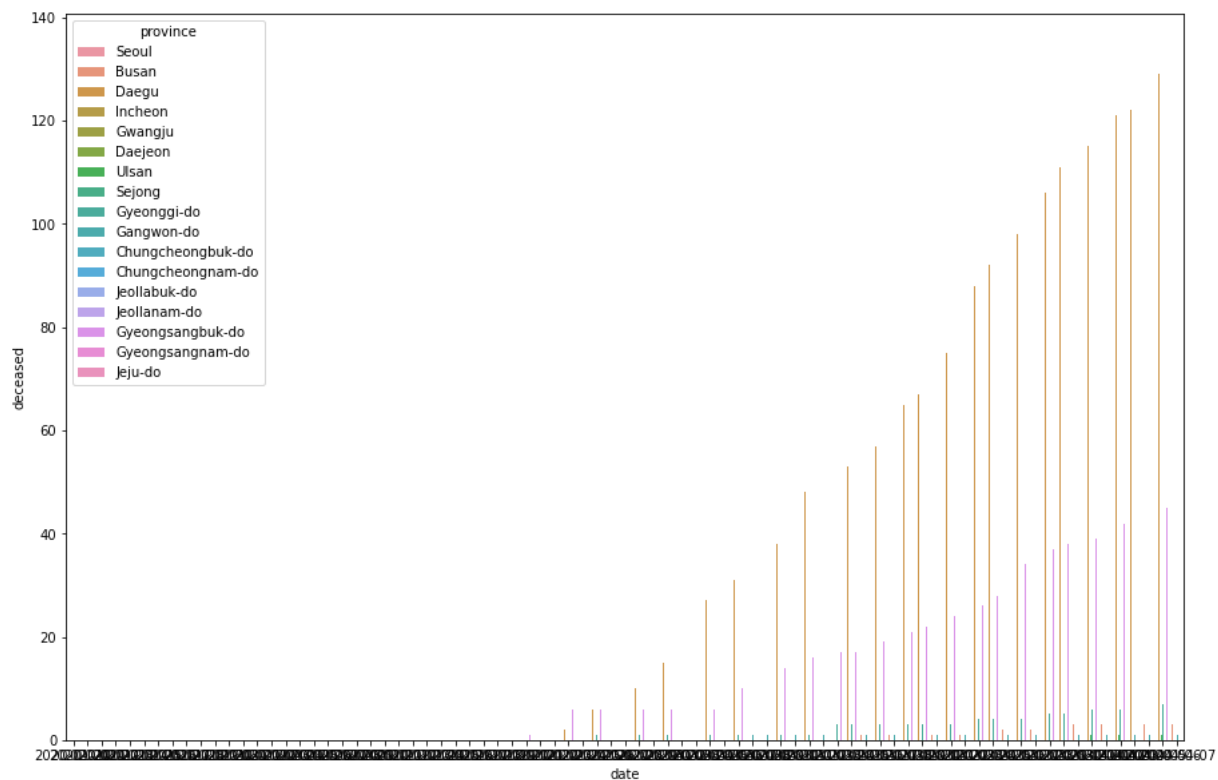
对released和province进行关联可视化处理

```
1 plt.figure(figsize=(15, 10))
2 ax=sn.barplot(x='date',y='released',hue = 'province',data=timeprovince)
```



对deceased和province进行关联可视化处理

```
1 plt.figure(figsize=(15, 10))
2 ax=sns.barplot(x='date',y='deceased',hue = 'province',data=timeprovince)
```



Task2 预测病人的恢复时间

首先，删除released_date为none的数据

```
1 list=patientdata[(patientdata['released_date']=='None').index.tolist()]
2 patientdata=patientdata.drop(list)
```

对confirmed_date和released_date列提取新特征，组成新的列

```

1 from datetime import datetime
2 patientdata['confirmed_date'].astype(str)
3 patientdata["confirm_year"] = patientdata.confirmed_date.apply(lambda x: x.split()[0].split("-")[0])
4 patientdata["confrim_month"] = patientdata.confirmed_date.apply(lambda x: datetime.strptime(x, "%Y-%m-%d").month)
5 patientdata["confrim_day"] = patientdata.confirmed_date.apply(lambda x: datetime.strptime(x, "%Y-%m-%d").day)
6 patientdata['released_date'].astype(str)
7 patientdata["released_year"] = patientdata.released_date.apply(lambda x: x.split()[0].split("-")[0])
8 patientdata["released_month"] = patientdata.released_date.apply(lambda x: datetime.strptime(x, "%Y-%m-%d").month)
9 patientdata["released_day"] = patientdata.released_date.apply(lambda x: datetime.strptime(x, "%Y-%m-%d").day)
10 patientdata
11

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	sex	birth_year	country	province	city	disease	infection_case	infected_by	contact_number	symptom_onset_date	...
0	1	1964.0	4	15	36	None	22	None	75.000000	2020-01-22	...
1	1	1987.0	4	15	85	None	22	None	31.000000	None	...
2	1	1964.0	4	15	83	None	18	2.002e+09	17.000000	None	...
3	1	1991.0	4	15	86	None	22	None	9.000000	2020-01-26	...
4	0	1992.0	4	15	112	None	18	1e+09	2.000000	None	...
...
3092	1	1971.0	4	9	39	None	19	None	18.908319	None	...
3119	1	1998.0	4	11	78	None	19	None	87.000000	None	...
3120	0	1998.0	4	11	78	None	19	None	84.000000	None	...
3121	0	1972.0	4	11	146	None	19	None	21.000000	None	...
3122	1	1974.0	4	11	78	None	19	None	74.000000	None	...

981 rows × 25 columns

计算confirm_date 和released_date 之间的日子，形成新的列recoverytime

定义函数计算recoverytime

```

1 def cal_Days(am,ad,fd):
2     if am==2 and fm==2:
3         return ad-fd
4     elif am==3 and fm==2:
5         return ad-fd+29
6     elif am==4 and fm==2:
7         return ad-fd+29+31
8     elif am==3 and fm==3:
9         return ad-fd
10    elif am==4 and fm==3:
11        return ad-fd+31
12    elif am==4 and fm==4:
13        return ad-fd
14    elif am==1 and fm==1:
15        return ad-fd
16    elif am==2 and fm==1:
17        return ad-fd+31
18    elif am==3 and fm==1:
19        return ad-fd+31+29
20    elif am==4 and fm==1:
21        return ad-fd+31+29+31

```

```

1 patientdata['recovery_time']=patientdata.apply(lambda x:
cal_Days(x['released_month'],x['released_day'],x['confrim_month'],x['confrim_day']),axis=1)

```

我们使用seaborn中的热力图来表明recovery_time和其他特征的关系，这里我们绘制recovery_time vs sex, birthyear, province, country, city, infection_case的热力图

```

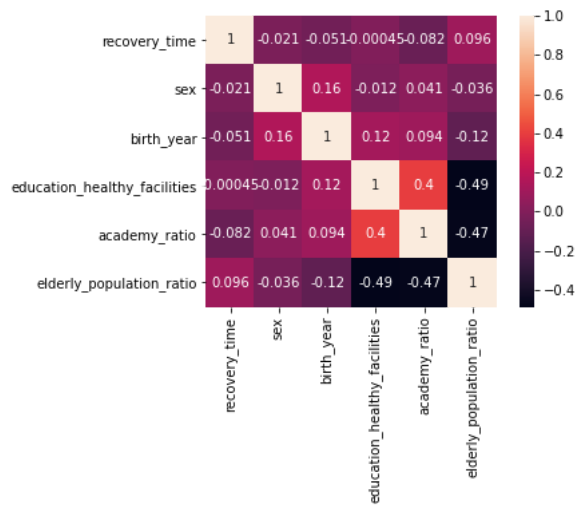
1 corrMat = patientdata[["recovery_time",
2 "sex","birth_year",'education_healthy_facilities','academy_ratio','elderly_population_ratio']].corr()
3 sn.heatmap(corrMat, annot=True, square=True)

```

```

1 <matplotlib.axes._subplots.AxesSubplot at 0x180d037d6d0>

```



我们可以看出academy_ratio和education_healthy_ratio之间有很高的相关性，我们关注的是recovery_time和其他各个特征之间的相关性
我们看出recovery_time和elderly_population_ratio正相关，和academy_ratio,education_healthy_facilities,birthyear

接下来通过这几种特征和线性回归预测病人的恢复时间，首先去除其他无用的特征

```

1 X=patientdata[["sex",'birth_year','education_healthy_facilities','academy_ratio','elderly_population_ratio']]
2 Ylabel=patientdata['recovery_time']
3 X.describe()

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	sex	birth_year	education_healthy_facilities	academy_ratio	elderly_population_ratio
count	981.000000	981.000000	981.000000	981.000000	981.000000
mean	0.405708	1979.910296	111.426096	1.497105	15.302691
std	0.491279	18.239453	56.889419	0.451158	4.859581
min	0.000000	1925.000000	0.000000	0.000000	0.000000
25%	0.000000	1966.000000	92.000000	1.340000	12.770000
50%	0.000000	1983.000000	102.000000	1.340000	16.180000
75%	1.000000	1995.000000	111.000000	1.760000	16.180000
max	1.000000	2020.000000	310.000000	4.180000	35.260000

训练集，测试集划分，比例为70，30

```

1 from sklearn.model_selection import train_test_split
2 train_x, test_x, train_y, test_y = train_test_split(X, Ylabel, test_size = 0.3, random_state=42)

```

```

1 # 构建线性回归模型
2 from sklearn.linear_model import LinearRegression
3 Model = LinearRegression()

```

```

1 X.isnull().sum()

```

```

1 sex 0
2 birth_year 0
3 education_healthy_facilities 0
4 academy_ratio 0
5 elderly_population_ratio 0
6 dtype: int64

```

```

1 # 用训练集训练模型
2 Model.fit(X= train_x, y = train_y)
3
4 # 模型预测
5 preds = Model.predict(X= test_x)

```

```

1 from sklearn.metrics import mean_squared_error
2 print("线性回归的标准差RMSE:", mean_squared_error(test_y, preds))

```

```

1 线性回归的标准差RMSE: 54.342400462634934

```

```

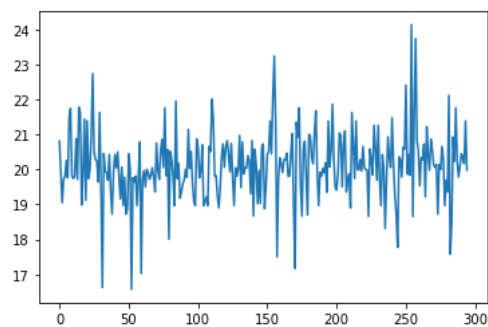
1 #结果可视化
2 plt.plot(preds)

```

```

1 [<matplotlib.lines.Line2D at 0x180badff460>]

```



和测试集的recovery_time相比较

```

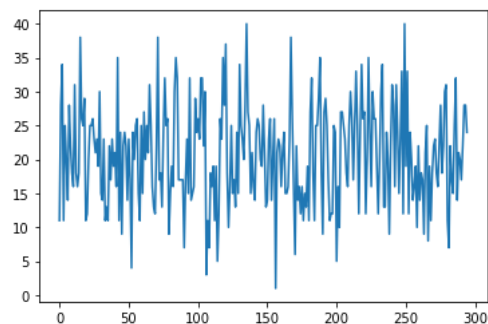
1 plt.plot(np.array(test_y))

```

```

1 [<matplotlib.lines.Line2D at 0x180ce945af0>]

```



绘制recovery_time vs birth_year education_healthy_facilities academy_ratio elderly_population_ratio的回归图

```

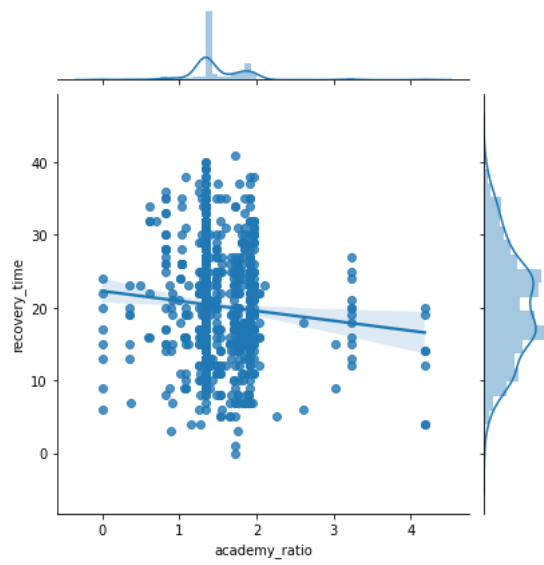
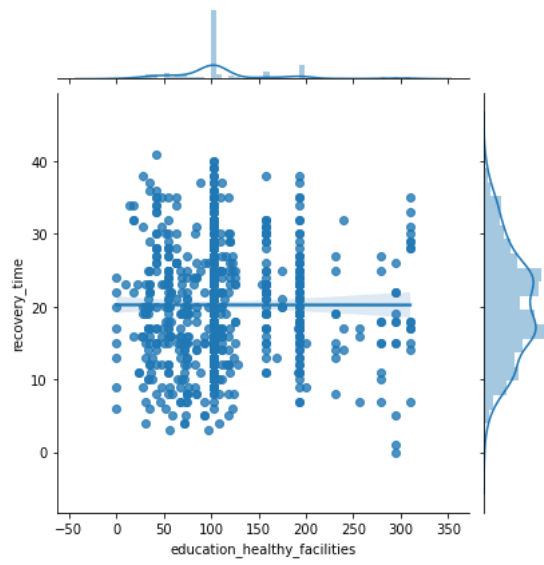
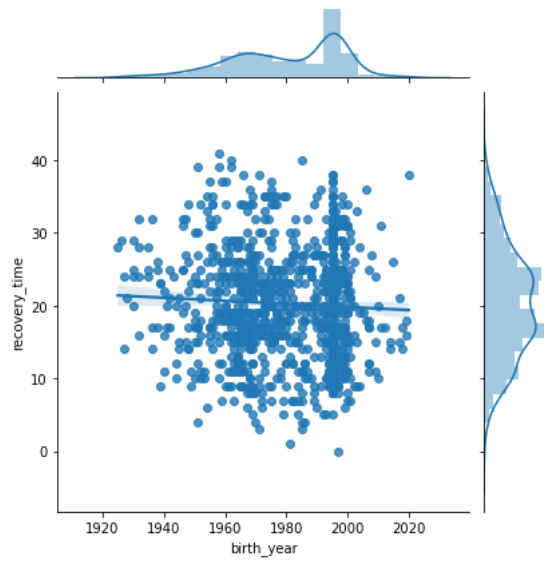
1 sn.jointplot(x="birth_year", y="recovery_time", data=patientdata, kind="reg")
2 sn.jointplot(x="education_healthy_facilities", y="recovery_time", data=patientdata, kind="reg")
3 sn.jointplot(x="academy_ratio", y="recovery_time", data=patientdata, kind="reg")
4 sn.jointplot(x="elderly_population_ratio", y="recovery_time", data=patientdata, kind="reg")

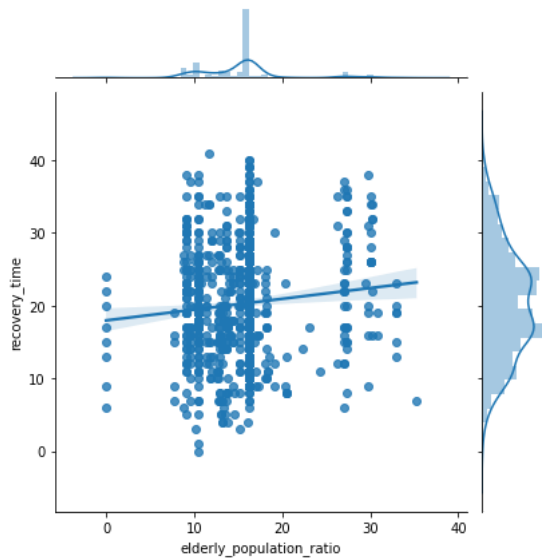
```

```

1 <seaborn.axisgrid.JointGrid at 0x180ce98c4f0>

```

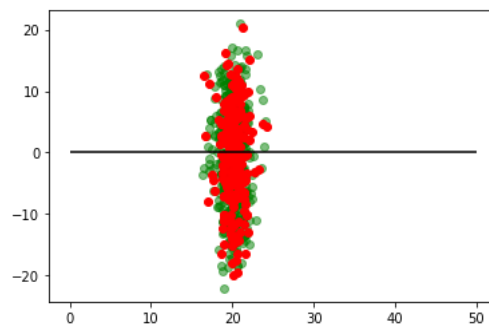




绘制残差图来诊断回归模型的效果，如果随机点分布在0的附近，就说明回归的效果好

```
1 plt.scatter(Model.predict(train_x), Model.predict(train_x) - train_y, c="g", alpha=0.5)
2 plt.scatter(Model.predict(test_x), Model.predict(test_x) - test_y, c="r")
3 plt.hlines(y=0, xmin=0, xmax=50)
4
```

```
1 <matplotlib.collections.LineCollection at 0x180cdcc2730>
```



task 3 城市的风险水平评估

1. 数据预处理

读入数据

```
1 Region=pd.read_csv("D:360Downloads\Region.csv")
2 TimeProvince=pd.read_csv("D:360Downloads\TimeProvince.csv")
3 PatientInfo=pd.read_csv("D:360Downloads\PatientInfo.csv")
4 Case=pd.read_csv("D:360Downloads\Case.csv")s
```

```
1 File "<ipython-input-86-e04ca505335d>", line 4
2 Case=pd.read_csv("D:360Downloads\Case.csv")s
3                                     ^
4 SyntaxError: invalid syntax
```

只获取最后一天的数据

```
1 list = TimeProvince[(TimeProvince['date']!='2020-04-07')].index.tolist()
2 TimeProvince = TimeProvince.drop(list)
3 TimeProvince.head(5)
```

```

1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }

```

	date	time	province	confirmed	released	deceased
1326	2020-04-07	0	Seoul	567	164	0
1327	2020-04-07	0	Busan	123	91	3
1328	2020-04-07	0	Daegu	6794	4918	134
1329	2020-04-07	0	Incheon	80	27	0
1330	2020-04-07	0	Gwangju	27	15	0

对 Region 文件进行预处理，删除经纬度、编码，添加确诊和死亡数

```

1 | # PatientCity = Region.dropcolumn(Region['latitude'].index.tolist())
2 |
3 | PatientCity = Region.drop(labels='latitude',axis=1, index=None, columns=None, inplace=False)
4 | PatientCity = PatientCity.drop(labels='longitude',axis=1, index=None, columns=None, inplace=False)
5 | PatientCity = PatientCity.drop(labels='code',axis=1, index=None, columns=None, inplace=False)
6 | for i in PatientCity.index:
7 |     if PatientCity.loc[i,'province'] == PatientCity.loc[i,'city']:
8 |         if (PatientCity.loc[i,'province'] == 'Sejong')or(PatientCity.loc[i,'province'] == 'Jeju-do')or(PatientCity.loc[i,'province'] ==
9 |         'Korea'):
10 |             continue
11 |         PatientCity = PatientCity.drop(labels=None,axis=0, index=[i], columns=None, inplace=False)
12 | PatientCity['oversea'] = 0
13 | PatientCity['group'] = 0
14 | PatientCity['confirmed'] = 0
15 | PatientCity['deceased'] = 0
16 | # PatientCity.loc[PatientCity.shape[0]] = [0 for n in range(PatientCity.shape[1])]
17 | # PatientCity.loc[244,'city'] = 'etc'
18 | PatientCity.head(5)

```

```

1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }

```

	province	city	elementary_school_count	kindergarten_count	university_count	academy_ratio	elderly_population_ratio	elderly_population_ratio
1	Seoul	Gangnam-gu	33	38	0	4.18	13.17	4.3
2	Seoul	Gangdong-gu	27	32	0	1.54	14.55	5.4
3	Seoul	Gangbuk-gu	14	21	0	0.67	19.49	8.5
4	Seoul	Gangseo-gu	36	56	1	1.17	14.39	5.7
5	Seoul	Gwanak-gu	22	33	1	0.89	15.12	4.9

创建map 用于匹配 PatientCity 里的行号

```

1 | map = {}
2 | for index,row in PatientCity.iterrows():
3 |     city = row['city']
4 |     map[city] = index

```

通过patientinfo获得每个城市的确诊死亡人数

```

1 | for index, row in PatientInfo.iterrows():
2 |     city = row["city"]
3 |     infect = row['infection_case']

```

```

4     if city is np.nan:
5         continue
6     if city == 'etc':
7         continue
8     ind = map[city]
9     if infect == 'overseas inflow':
10         PatientCity.loc[ind,'oversea'] = PatientCity.loc[ind,'oversea'] + 1
11     PatientCity.loc[ind,'confirmed'] = PatientCity.loc[ind,'confirmed'] + 1
12     if row["deceased_date"] is not np.nan:
13         PatientCity.loc[ind,'deceased'] = PatientCity.iloc[ind]['deceased'] + 1
14
15
16 for index, row in Case.iterrows():
17     city = row["city"]
18     Group = row['group']
19     num = row['confirmed'];
20     if city is np.nan:
21         continue
22     if city == 'etc':
23         continue
24     if city == 'from other city':
25         continue
26     if city == '-':
27         continue
28     if city == 'Jin-gu':
29         city = 'Busanjin-gu'
30     ind = map[city]
31     if Group == True:
32         PatientCity.loc[ind,'group'] = PatientCity.loc[ind,'group'] + num;
33 PatientCity.head()

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	province	city	elementary_school_count	kindergarten_count	university_count	academy_ratio	elderly_population_ratio	elderly_alone_ratio
1	Seoul	Gangnam-gu	33	38	0	4.18	13.17	4.3
2	Seoul	Gangdong-gu	27	32	0	1.54	14.55	5.4
3	Seoul	Gangbuk-gu	14	21	0	0.67	19.49	8.5
4	Seoul	Gangseo-gu	36	56	1	1.17	14.39	5.7
5	Seoul	Gwanak-gu	22	33	1	0.89	15.12	4.9

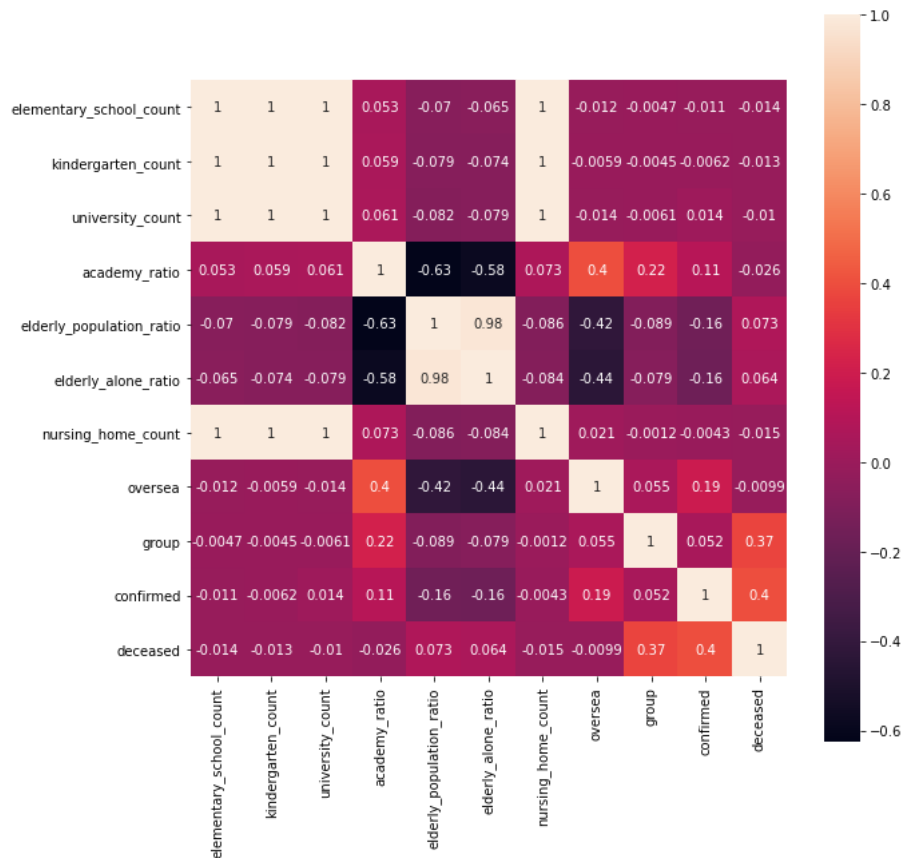
用热力图判断相关性

```

1 plt.figure(figsize=(10,10))
2 corrMat = PatientCity[["elementary_school_count",
3     "kindergarten_count",
4     "university_count",
5     'academy_ratio',
6     'elderly_population_ratio',
7     'elderly_alone_ratio',
8     'nursing_home_count',
9     'oversea',
10    'group',
11    'confirmed',
12    'deceased']].corr()
13 sn.heatmap(corrMat, annot=True, square=True)

```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x180d10e3dc0>
```



发现:

"elementary_school_count", "kindergarten_count", "university_count"与'nursing_home_count' 有较高的相关性

'elderly_population_ratio', 'elderly_alone_ratio' 之间有较高的相关性。

因此

将 "elementary_school_count", "kindergarten_count", "university_count"与'nursing_home_count' 四个只保留 nursing_home_count;

'elderly_population_ratio', 'elderly_alone_ratio' 只保留 elderly_population_ratio

```
1 PatientCityWithProvince = PatientCity
2 PatientCity = PatientCity.drop(labels='province',axis=1, index=None, columns=None, inplace=False)
3 PatientCity.head(10)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	city	elementary_school_count	kindergarten_count	university_count	academy_ratio	elderly_population_ratio	elderly_alone_
1	Gangnam-gu	33	38	0	4.18	13.17	4.3
2	Gangdong-gu	27	32	0	1.54	14.55	5.4
3	Gangbuk-gu	14	21	0	0.67	19.49	8.5
4	Gangseo-gu	36	56	1	1.17	14.39	5.7
5	Gwanak-gu	22	33	1	0.89	15.12	4.9
6	Gwangjin-gu	22	33	3	1.16	13.75	4.8
7	Guro-gu	26	34	3	1.00	16.21	5.7
8	Geumcheon-gu	18	19	0	0.96	16.15	6.7
9	Nowon-gu	42	66	6	1.39	15.40	7.4
10	Dobong-gu	23	26	1	0.95	17.89	7.2

```

1 PatientCityWithProvince = PatientCityWithProvince.drop(labels='elementary_school_count',axis=1, index=None, columns=None, inplace=False)
2 PatientCityWithProvince = PatientCityWithProvince.drop(labels='kindergarten_count',axis=1, index=None, columns=None, inplace=False)
3 PatientCityWithProvince = PatientCityWithProvince.drop(labels='university_count',axis=1, index=None, columns=None, inplace=False)
4 PatientCityWithProvince = PatientCityWithProvince.drop(labels='elderly_alone_ratio',axis=1, index=None, columns=None, inplace=False)
5 PatientCityWithProvince.head(10)

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	province	city	academy_ratio	elderly_population_ratio	nursing_home_count	oversea	group	confirmed	deceased
1	Seoul	Gangnam-gu	4.18	13.17	3088	31	0	53	0
2	Seoul	Gangdong-gu	1.54	14.55	1023	8	0	13	0
3	Seoul	Gangbuk-gu	0.67	19.49	628	1	0	5	0
4	Seoul	Gangseo-gu	1.17	14.39	1080	0	0	0	0
5	Seoul	Gwanak-gu	0.89	15.12	909	10	0	41	0
6	Seoul	Gwangjin-gu	1.16	13.75	723	6	0	7	0
7	Seoul	Guro-gu	1.00	16.21	741	3	138	34	0
8	Seoul	Geumcheon-gu	0.96	16.15	475	1	0	13	0
9	Seoul	Nowon-gu	1.39	15.40	952	5	0	23	0
10	Seoul	Dobong-gu	0.95	17.89	485	4	0	7	0

2.数据处理：

导入数据

```

1 from sklearn.pipeline import make_pipeline
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.cluster import KMeans
4 import pandas as pd
5
6 riskLevels = PatientCity.drop(labels='elementary_school_count',axis=1, index=None, columns=None, inplace=False)
7 riskLevels = riskLevels.drop(labels='kindergarten_count',axis=1, index=None, columns=None, inplace=False)
8 riskLevels = riskLevels.drop(labels='university_count',axis=1, index=None, columns=None, inplace=False)
9 riskLevels = riskLevels.drop(labels='elderly_alone_ratio',axis=1, index=None, columns=None, inplace=False)
10
11 df=riskLevels
12 grain_variety=df['city']
13 del df['city']
14 df.head()

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	academy_ratio	elderly_population_ratio	nursing_home_count	oversea	group	confirmed	deceased
1	4.18	13.17	3088	31	0	53	0
2	1.54	14.55	1023	8	0	13	0
3	0.67	19.49	628	1	0	5	0
4	1.17	14.39	1080	0	0	0	0
5	0.89	15.12	909	10	0	41	0

数据标准化和kmeans分析

```

1 samples=df.values
2 #标准化
3 scaler=StandardScaler()
4 kmeans=KMeans(n_clusters=3)
5 pipeline=make_pipeline(scaler,kmeans)
6 pipeline.fit(samples)#训练模型
7 labels=pipeline.predict(samples)#预测
8 df=pd.DataFrame({'labels':labels,'city':grain_variety})
9 ct=pd.crosstab(df['labels'],df['city'])
10 ct

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

city	Andong-si	Ansan-si	Anseong-si	Anyang-si	Asan-si	Boeun-gun	Bonghwa-gun	Boryeong-si	Boseong-gun	Buan-gun	...	Yeongju-si	Yeongw g
labels													
0	0	0	0	0	0	0	1	0	0	0	...	0	0
1	0	0	0	0	0	1	0	1	1	1	...	1	1
2	1	1	1	1	1	0	0	0	0	0	...	0	0

3 rows × 207 columns

```

1 # 获取某一行的方法:
2 def getRow(a,data):
3     return data[a:a+1]

```

选择这些列表里三个不同 label 的 city 查看

结果如下:

label==2

```

1 for index,row in PatientCity.iterrows():
2     if row['city'] == 'Ansan-si':
3         print(row)

```

```

1 city          Ansan-si
2 elementary_school_count    54
3 kindergarten_count        94
4 university_count           4
5 academy_ratio             1.49
6 elderly_population_ratio   10.35
7 elderly_alone_ratio        4.6
8 nursing_home_count        1024
9 oversea                   3
10 group                    0
11 confirmed                16
12 deceased                 0
13 Name: 97, dtype: object

```

label==1

```

1 for index,row in PatientCity.iterrows():
2     if row['city'] == 'Nam-gu':
3         print(row)

```

```

1 city          Nam-gu
2 elementary_school_count    21
3 kindergarten_count        27
4 university_count           4
5 academy_ratio             1.24
6 elderly_population_ratio   19.13
7 elderly_alone_ratio        7.9
8 nursing_home_count        475
9 oversea                   0

```

```

10 | group                0
11 | confirmed            0
12 | deceased            0
13 | Name: 30, dtype: object
14 | city                Nam-gu
15 | elementary_school_count  11
16 | kindergarten_count    15
17 | university_count      2
18 | academy_ratio         0.85
19 | elderly_population_ratio 22.49
20 | elderly_alone_ratio    10.4
21 | nursing_home_count    345
22 | oversea              0
23 | group                0
24 | confirmed            0
25 | deceased            0
26 | Name: 44, dtype: object
27 | city                Nam-gu
28 | elementary_school_count  23
29 | kindergarten_count    44
30 | university_count      4
31 | academy_ratio         2.63
32 | elderly_population_ratio 16.76
33 | elderly_alone_ratio    7.5
34 | nursing_home_count    427
35 | oversea              0
36 | group                0
37 | confirmed            0
38 | deceased            0
39 | Name: 54, dtype: object
40 | city                Nam-gu
41 | elementary_school_count  30
42 | kindergarten_count    44
43 | university_count      1
44 | academy_ratio         3.23
45 | elderly_population_ratio 11.42
46 | elderly_alone_ratio    4.8
47 | nursing_home_count    765
48 | oversea              5
49 | group                4508
50 | confirmed            28
51 | deceased            1
52 | Name: 76, dtype: object

```

label==1

```

1 | for index,row in PatientCity.iterrows():
2 |     if row['city'] == 'Bonghwa-gun':
3 |         print(row)

```

```

1 | city                Bonghwa-gun
2 | elementary_school_count  14
3 | kindergarten_count    17
4 | university_count      0
5 | academy_ratio         0.37
6 | elderly_population_ratio 35.26
7 | elderly_alone_ratio    20
8 | nursing_home_count    47
9 | oversea              0
10 | group                68
11 | confirmed            71
12 | deceased            1
13 | Name: 207, dtype: object

```

从上边的例子里边 可以看出来,

label == 2 时是风险最低的情况;

label == 0 时是风险中度的情况

label == 1 是的风险最高的情况

高风险地区:

```

1 | high = getRow(1,ct)
2 | city = []
3 | cols=[x for i,x in enumerate(high.columns) if high.iat[0,i]==0]
4 | high=high.drop(cols,axis=1) #利用drop方法将含有特定数值的列删除
5 | print("高风险地区:")
6 |
7 | for i,x in enumerate(high.columns):
8 |     city.append(x)
9 | output = PatientCitywithProvince
10 | for index,row in output.iterrows():
11 |     if row['city'] not in city:
12 |         output = output.drop(index,axis=0)

```

```

13 for index,row in output.iterrows():
14     if row['confirmed'] > 20:
15         print(row['province'],row['city'])

```

```

1 高风险地区:
2  Daejeon Seo-gu
3  Ulsan Nam-gu
4  Ulsan Jung-gu
5  Gyeongsangbuk-do Gimcheon-si
6  Gyeongsangbuk-do Uiseong-gun

```

中风险地区:

```

1  mid = getRow(0,ct)
2  city = []
3  cols=[x for i,x in enumerate(mid.columns) if mid.iat[0,i]==0]
4  mid=mid.drop(cols,axis=1) #利用drop方法将含有特定数值的列删除
5  print("中风险地区:")
6  for i,x in enumerate(mid.columns):
7      city.append(x)
8  output = PatientCitywithProvince
9  for index,row in output.iterrows():
10     if row['city'] not in city:
11         output = output.drop(index,axis=0)
12 for index,row in output.iterrows():
13     print(row['province'],row['city'])

```

```

1  中风险地区:
2  Busan Nam-gu
3  Daegu Nam-gu
4  Gwangju Nam-gu
5  Ulsan Nam-gu
6  Gyeonggi-do Namyangju-si
7  Gangwon-do Taebaek-si
8  Gyeongsangbuk-do Gyeongsan-si
9  Gyeongsangbuk-do Bonghwa-gun
10 Gyeongsangbuk-do Seongju-gun
11 Gyeongsangbuk-do Yeongcheon-si
12 Gyeongsangbuk-do Cheongdo-gun

```

低风险地区:

```

1  low = getRow(2,ct)
2  city = []
3  cols=[x for i,x in enumerate(low.columns) if low.iat[0,i]==0]
4  low=low.drop(cols,axis=1) #利用drop方法将含有特定数值的列删除
5  print("低风险地区:")
6  for i,x in enumerate(low.columns):
7      city.append(x)
8  output = PatientCitywithProvince
9  for index,row in output.iterrows():
10     if row['city'] not in city:
11         output = output.drop(index,axis=0)
12 for index,row in output.iterrows():
13     print(row['province'],row['city'])

```

```

1  低风险地区:
2  Seoul Gangnam-gu
3  Seoul Gangdong-gu
4  Seoul Gangseo-gu
5  Seoul Gwanak-gu
6  Seoul Gwangjin-gu
7  Seoul Guro-gu
8  Seoul Geumcheon-gu
9  Seoul Nowon-gu
10 Seoul Dobong-gu
11 Seoul Dongdaemun-gu
12 Seoul Dongjak-gu
13 Seoul Mapo-gu
14 Seoul Seodaemun-gu
15 Seoul Seocho-gu
16 Seoul Seongdong-gu
17 Seoul Seongbuk-gu
18 Seoul Songpa-gu
19 Seoul Yangcheon-gu
20 Seoul Yeongdeungpo-gu
21 Seoul Yongsan-gu
22 Seoul Eunpyeong-gu
23 Seoul Jongno-gu
24 Seoul Jung-gu
25 Seoul Jungnang-gu
26 Busan Gangseo-gu
27 Busan Geumjeong-gu

```


28	Busan Gijang-gun
29	Busan Nam-gu
30	Busan Dong-gu
31	Busan Dongnae-gu
32	Busan Busanjin-gu
33	Busan Buk-gu
34	Busan Saha-gu
35	Busan Seo-gu
36	Busan Suyeong-gu
37	Busan Yeonje-gu
38	Busan Jung-gu
39	Busan Haeundae-gu
40	Daegu Nam-gu
41	Daegu Dalseo-gu
42	Daegu Dalseong-gun
43	Daegu Dong-gu
44	Daegu Buk-gu
45	Daegu Seo-gu
46	Daegu Suseong-gu
47	Daegu Jung-gu
48	Gwangju Gwangsan-gu
49	Gwangju Nam-gu
50	Gwangju Dong-gu
51	Gwangju Buk-gu
52	Gwangju Seo-gu
53	Incheon Gyeyang-gu
54	Incheon Michuhol-gu
55	Incheon Namdong-gu
56	Incheon Dong-gu
57	Incheon Bupyeong-gu
58	Incheon Seo-gu
59	Incheon Yeonsu-gu
60	Incheon Jung-gu
61	Daejeon Daedeok-gu
62	Daejeon Dong-gu
63	Daejeon Seo-gu
64	Daejeon Yuseong-gu
65	Daejeon Jung-gu
66	Ulsan Nam-gu
67	Ulsan Dong-gu
68	Ulsan Buk-gu
69	Ulsan Ulju-gun
70	Ulsan Jung-gu
71	Sejong Sejong
72	Gyeonggi-do Goyang-si
73	Gyeonggi-do Gwacheon-si
74	Gyeonggi-do Gwangmyeong-si
75	Gyeonggi-do Gwangju-si
76	Gyeonggi-do Guri-si
77	Gyeonggi-do Gunpo-si
78	Gyeonggi-do Gimpo-si
79	Gyeonggi-do Bucheon-si
80	Gyeonggi-do Seongnam-si
81	Gyeonggi-do Suwon-si
82	Gyeonggi-do Siheung-si
83	Gyeonggi-do Ansan-si
84	Gyeonggi-do Anseong-si
85	Gyeonggi-do Anyang-si
86	Gyeonggi-do Yangju-si
87	Gyeonggi-do Osan-si
88	Gyeonggi-do Yongin-si
89	Gyeonggi-do Uiwang-si
90	Gyeonggi-do Uijeongbu-si
91	Gyeonggi-do Icheon-si
92	Gyeonggi-do Paju-si
93	Gyeonggi-do Pyeongtaek-si
94	Gyeonggi-do Hanam-si
95	Gyeonggi-do Hwaseong-si
96	Gangwon-do Gangneung-si
97	Gangwon-do Donghae-si
98	Gangwon-do Sokcho-si
99	Gangwon-do Wonju-si
100	Gangwon-do Chuncheon-si
101	Chungcheongbuk-do Jeungpyeong-gun
102	Chungcheongbuk-do Jincheon-gun
103	Chungcheongbuk-do Cheongju-si
104	Chungcheongbuk-do Chungju-si
105	Chungcheongnam-do Gyeryong-si
106	Chungcheongnam-do Dangjin-si
107	Chungcheongnam-do Seosan-si
108	Chungcheongnam-do Asan-si
109	Chungcheongnam-do Cheonan-si
110	Jeollabuk-do Gunsan-si
111	Jeollabuk-do Iksan-si
112	Jeollabuk-do Jeonju-si
113	Jeollanam-do Gwangyang-si
114	Jeollanam-do Mokpo-si

```

115 Jeollanam-do Muan-gun
116 Jeollanam-do Suncheon-si
117 Jeollanam-do Yeosu-si
118 Gyeongsangbuk-do Gyeongju-si
119 Gyeongsangbuk-do Gumi-si
120 Gyeongsangbuk-do Andong-si
121 Gyeongsangbuk-do Chilgok-gun
122 Gyeongsangbuk-do Pohang-si
123 Gyeongsangnam-do Geoje-si
124 Gyeongsangnam-do Gimhae-si
125 Gyeongsangnam-do Sacheon-si
126 Gyeongsangnam-do Yangsan-si
127 Gyeongsangnam-do Jinju-si
128 Gyeongsangnam-do Changwon-si
129 Gyeongsangnam-do Tongyeong-si
130 Jeju-do Jeju-do
131 Korea Korea

```

Task4 预测感染人数

我们先来绘制疫情确诊人数的走势图

```

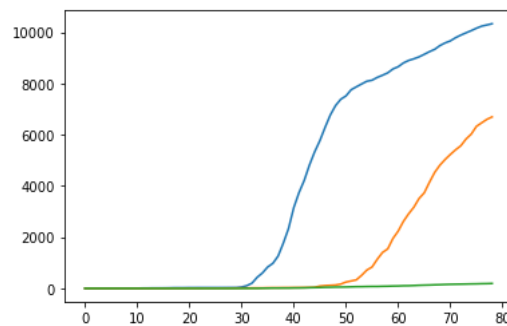
1 plt.plot(time['confirmed'],label='confirmed')
2 plt.plot(time['released'],label='released')
3 plt.plot(time['deceased'],label='deceased')
4

```

```

1 [

```



我们考虑用逻辑斯蒂函数来预测

```

1 #定义逻辑斯蒂函数
2 def logistic_function(p,t):
3     k,b=p
4     a=0.1
5
6     value=np.exp(-a*(t-b))
7     return k/(1+value)
8 #定义误差函数
9 def error_f(p,t,y):
10     return logistic_function(p,t)-y
11 time.info()
12
13

```

```

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 79 entries, 0 to 78
3 Data columns (total 7 columns):
4 #   Column      Non-Null Count  Dtype
5 ---  ---
6 0    date        79 non-null    object
7 1    time        79 non-null    int64
8 2    test        79 non-null    int64
9 3    negative    79 non-null    int64
10 4    confirmed   79 non-null    int64
11 5    released    79 non-null    int64
12 6    deceased    79 non-null    int64
13 dtypes: int64(6), object(1)
14 memory usage: 4.4+ KB

```

```

1 from scipy.optimize import leastsq
2 #参数初始值

```

```

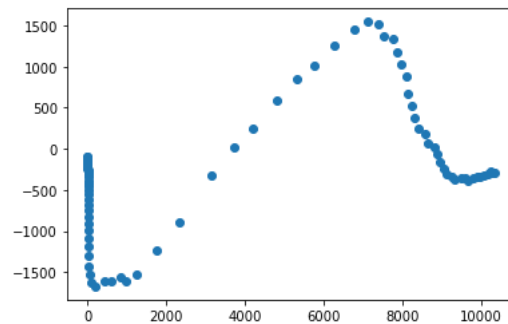
3 logistic_p0=[10000,20]
4 t=np.array([i+1 for i in range(79)])
5 y=time['confirmed'].values
6 #利用最小二乘法求解参数
7 logistic_paramtr=leastsq(error_f,logistic_p0,args=(t,y))
8 p=logistic_paramtr[0]
9 #利用我们定义的逻辑斯蒂函数进行预测
10 predict_data=logistic_function(p,t)
11 predict_data
12 #预测的误差
13 error=y-predict_data
14 error
15 #绘制误差的散点图
16 plt.scatter(y,error)
17

```

```

1 <matplotlib.collections.PathCollection at 0x180ce4ef3a0>

```



```

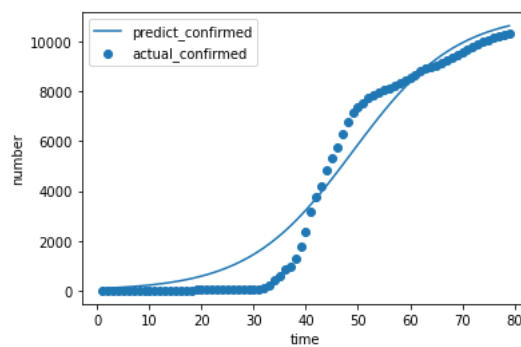
1 #绘图
2 plt.scatter(t,y,label='actual_confirmed')
3 plt.plot(t,predict_data,label='predict_confirmed')
4 plt.xlabel('time')
5 plt.ylabel('number')
6 plt.legend(loc='best')
7

```

```

1 <matplotlib.legend.Legend at 0x180cd9ac670>

```

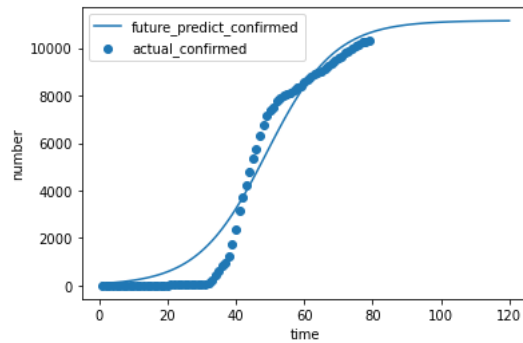


```

1 #预测未来疫情的走势
2 future_t=[i+1 for i in range(120)]
3 future_predict=logistic_function(p,future_t)
4 future_predict
5 #绘图
6 #绘图
7 plt.scatter(t,y,label='actual_confirmed')
8 plt.plot(future_t,future_predict,label='future_predict_confirmed')
9 plt.xlabel('time')
10 plt.ylabel('number')
11 plt.legend(loc='best')

```

```
1 | <matplotlib.legend.Legend at 0x180ce58bc40>
```



预测接下来韩国的确诊人数

```
1 | for i in range(81,119):
2 |     print(future_predict[i])
3 |
```

```
1 | 10760.21409500655
2 | 10796.652288496583
3 | 10829.836260428989
4 | 10860.03864093677
5 | 10887.512426223213
6 | 10912.491803741863
7 | 10935.19306332525
8 | 10955.815562195587
9 | 10974.542717779608
10 | 10991.54300745268
11 | 11006.970958800215
12 | 11020.968117766812
13 | 11033.663985233312
14 | 11045.176915189058
15 | 11055.614969819364
16 | 11065.076728571732
17 | 11073.652049658514
18 | 11081.422783552653
19 | 11088.463438884637
20 | 11094.841801795546
21 | 11100.619510278748
22 | 11105.852585383465
23 | 11110.591921383004
24 | 11114.883737151555
25 | 11118.76999106493
26 | 11122.288761757494
27 | 11125.474597043065
28 | 11128.358833251825
29 | 11130.969887156922
30 | 11133.333522570134
31 | 11135.473093581173
32 | 11137.409766304334
33 | 11139.16272088227
34 | 11140.749335382636
35 | 11142.1853531109
36 | 11143.485034753186
37 | 11144.661296657818
38 | 11145.725836463826
```

因此，预测韩国的疫情将于开始后的第120天结束，结束时的累计确诊量到达11145例

从经纬度找到感染的聚集点

```
1 | patientdata.head()
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

	sex	birth_year	country	province	city	disease	infection_case	infected_by	contact_number	symptom_onset_date	...	
0	1	1964.0	4	15	36	None	22	None	75.0	2020-01-22	...	1
1	1	1987.0	4	15	85	None	22	None	31.0	None	...	1
2	1	1964.0	4	15	83	None	18	2.002e+09	17.0	None	...	1
3	1	1991.0	4	15	86	None	22	None	9.0	2020-01-26	...	1
4	0	1992.0	4	15	112	None	18	1e+09	2.0	None	...	1

5 rows × 26 columns

```
1 newpatientdata=patientdata[['latitude','longitude']]
2 newpatientdata.head()
```

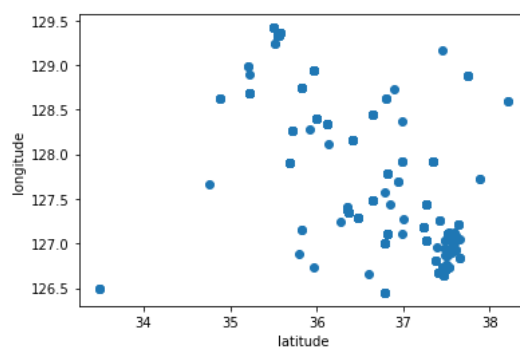
```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	latitude	longitude
0	35.212424	128.980680
1	37.606832	127.092656
2	37.572999	126.979189
3	37.566283	126.901644
4	37.589562	127.016700

用kmeans对该数据进行聚类，找到各个聚集的中心点

```
1 newpatientdata=newpatientdata[newpatientdata['latitude']!=0]
2
3 plt.scatter(newpatientdata['latitude'],newpatientdata['longitude'])
4 plt.xlabel('latitude')
5 plt.ylabel('longitude')
```

```
1 Text(0, 0.5, 'longitude')
```



```
1 from sklearn.cluster import KMeans
2 #构造聚类器
3 estimator=KMeans(n_clusters=5)
4 estimator.fit(newpatientdata)
5 label_predict=estimator.labels_
```

```

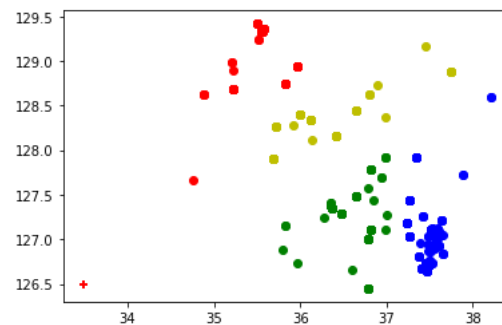
1 x0=newpatientdata[label_predict==0]
2 x1=newpatientdata[label_predict==1]
3 x2=newpatientdata[label_predict==2]
4 x3=newpatientdata[label_predict==3]
5 x4=newpatientdata[label_predict==4]
6 plt.scatter(x0['latitude'],x0['longitude'],c='r',label='x0')
7 plt.scatter(x1['latitude'],x1['longitude'],c='g',label='x1')
8 plt.scatter(x2['latitude'],x2['longitude'],c='b',label='x2')
9 plt.scatter(x3['latitude'],x3['longitude'],c='y',label='x3')
10 plt.scatter(x4['latitude'],x4['longitude'],c='r',marker='+',label='x4')

```

```

1 <matplotlib.collections.PathCollection at 0x180ce286c70>

```



task 判断当日确诊人数和当日的天气之间的关系，我们先选择韩国确诊人数最多的Gyeongsangbuk-do进行分析

```

1 timeprovince=pd.read_csv("D:360Downloads\TimeProvince.csv")
2 time_gye=timeprovince[timeprovince['province']=='Gyeongsangbuk-do']
3
4 weather=pd.read_csv("D:360Downloads\weather.csv")
5 weather_gye=weather[weather['province']=='Gyeongsangbuk-do']
6 #定义函数，将每天Gyeongsangbuk-do的天气数据导入，
7 def getweather(dataframe,dR):
8
9     dataframe['avg_temp']=0
10    dataframe['min_temp']=0
11    dataframe['max_temp']=0
12    dataframe['max_wind_speed']=0
13    dataframe['most_wind_direction']=0
14    dataframe['avg_relative_humidity']=0
15    for i in dataframe.index:
16        for j in dR.index:
17
18            if dataframe.loc[i,'date']==dR.loc[j,'date']:
19                dataframe.loc[i,'avg_temp']=dR.loc[j,'avg_temp']
20                dataframe.loc[i,'max_temp']=dR.loc[j,'max_temp']
21                dataframe.loc[i,'min_temp']=dR.loc[j,'min_temp']
22                dataframe.loc[i,'max_wind_speed']=dR.loc[j,'max_wind_speed']
23                dataframe.loc[i,'most_wind_direction']=dR.loc[j,'most_wind_direction']
24                dataframe.loc[i,'avg_relative_humidity']=dR.loc[j,'avg_relative_humidity']
25
26    return dataframe
27 getweather(time_gye,weather_gye)

```

```

1 <ipython-input-105-42a2d71124d3>:9: SettingWithCopyWarning:
2 A value is trying to be set on a copy of a slice from a DataFrame.
3 Try using .loc[row_indexer,col_indexer] = value instead
4
5 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
6     dataframe['avg_temp']=0
7 <ipython-input-105-42a2d71124d3>:10: SettingWithCopyWarning:
8 A value is trying to be set on a copy of a slice from a DataFrame.
9 Try using .loc[row_indexer,col_indexer] = value instead
10
11 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
12     dataframe['min_temp']=0
13 <ipython-input-105-42a2d71124d3>:11: SettingWithCopyWarning:
14 A value is trying to be set on a copy of a slice from a DataFrame.
15 Try using .loc[row_indexer,col_indexer] = value instead
16

```

```

17 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
18     dataframe['max_temp']=0
19 <ipython-input-105-42a2d71124d3>:12: SettingWithCopyWarning:
20 A value is trying to be set on a copy of a slice from a DataFrame.
21 Try using .loc[row_indexer,col_indexer] = value instead
22
23 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
24     dataframe['max_wind_speed']=0
25 <ipython-input-105-42a2d71124d3>:13: SettingWithCopyWarning:
26 A value is trying to be set on a copy of a slice from a DataFrame.
27 Try using .loc[row_indexer,col_indexer] = value instead
28
29 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
30     dataframe['most_wind_direction']=0
31 <ipython-input-105-42a2d71124d3>:14: SettingWithCopyWarning:
32 A value is trying to be set on a copy of a slice from a DataFrame.
33 Try using .loc[row_indexer,col_indexer] = value instead
34
35 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
36     dataframe['avg_relative_humidity']=0
37 d:\python\lib\site-packages\pandas\core\indexing.py:966: SettingWithCopyWarning:
38 A value is trying to be set on a copy of a slice from a DataFrame.
39 Try using .loc[row_indexer,col_indexer] = value instead
40
41 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
42     self.obj[item] = s

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	date	time	province	confirmed	released	deceased	avg_temp	min_temp	max_temp	max_wind_speed	most_w
14	2020-01-20	16	Gyeongsangbuk-do	0	0	0	3.6	-0.4	6.9	7.5	270.0
31	2020-01-21	16	Gyeongsangbuk-do	0	0	0	2.0	-1.9	7.5	4.8	180.0
48	2020-01-22	16	Gyeongsangbuk-do	0	0	0	5.2	-0.2	8.4	3.4	180.0
65	2020-01-23	16	Gyeongsangbuk-do	0	0	0	6.9	4.2	11.3	3.4	160.0
82	2020-01-24	16	Gyeongsangbuk-do	0	0	0	5.8	1.9	11.9	8.7	20.0
...
1272	2020-04-03	0	Gyeongsangbuk-do	1309	865	42	14.4	8.6	20.5	8.8	290.0
1289	2020-04-04	0	Gyeongsangbuk-do	1310	881	44	9.6	3.3	16.6	10.8	50.0
1306	2020-04-05	0	Gyeongsangbuk-do	1314	899	44	6.0	0.0	14.0	6.0	180.0
1323	2020-04-06	0	Gyeongsangbuk-do	1316	922	45	10.8	0.9	17.3	5.9	180.0
1340	2020-04-07	0	Gyeongsangbuk-do	1317	934	46	0.0	0.0	0.0	0.0	0.0

79 rows × 12 columns

```

1 #再去掉confirmed为0的数据，只查看有感染人数的数据
2 time_gye=time_gye[time_gye['confirmed']!=0]
3 new_time_gye=time_gye
4 new_time_gye.head()
5 for i in time_gye.index:
6     if i>524:
7         new_time_gye.loc[i,'confirmed']=time_gye.loc[i,'confirmed']-time_gye.loc[i-17,'confirmed']
8 new_time_gye.head()

```

```

1 d:\python\lib\site-packages\pandas\core\indexing.py:966: SettingWithCopyWarning:
2 A value is trying to be set on a copy of a slice from a DataFrame.
3 Try using .loc[row_indexer,col_indexer] = value instead
4
5 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
6 self.obj[item] = s

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	date	time	province	confirmed	released	deceased	avg_temp	min_temp	max_temp	max_wind_speed	most_wi
524	2020-02-19	16	Gyeongsangbuk-do	2	1	0	3.8	-1.1	11.2	10.0	180.0
541	2020-02-20	16	Gyeongsangbuk-do	23	1	1	5.8	1.0	12.2	4.6	180.0
558	2020-02-21	16	Gyeongsangbuk-do	5	1	1	8.3	-0.2	15.2	5.7	180.0
575	2020-02-22	16	Gyeongsangbuk-do	140	1	3	7.9	2.7	11.1	9.9	270.0
592	2020-02-23	16	Gyeongsangbuk-do	30	1	3	5.2	-0.5	10.5	11.5	270.0

用热力图研究确诊数量和天气因素的相关性

```

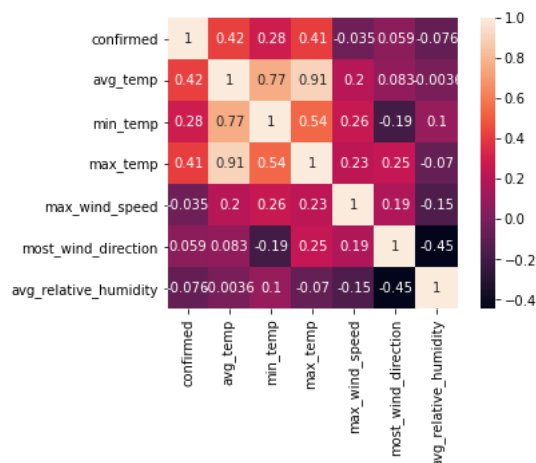
1 corrMat = new_time_gye[["confirmed",
2 "avg_temp", "min_temp", "max_temp", "max_wind_speed", "most_wind_direction", "avg_relative_humidity"]].corr()
3 sn.heatmap(corrMat, annot=True, square=True)

```

```

1 <matplotlib.axes._subplots.AxesSubplot at 0x180ce79d190>

```



可以看出确诊数量和avg_temp,min_temp,max_temp, 湿度, 风速正相关。

