

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm

import warnings
warnings.filterwarnings("ignore")
```

## Data Outline

```
In [2]: data=pd.read_excel("variable in TAZ.xlsx",index_col=0)
data.head()
```

```
Out[2]:
```

	Open Area	Traffic Area	Residential Area	Commercial Area	POI Entropy	POI Richness	POI Simpson	POI Gini Coefficient	Road Density	Population Density	Distance to Covid Breakout Location
FID											
0	71.845258	4.294096	13.019560	11.041090	5.091128	9	4.198463	0.639767	9.735901	5.890306	14.6319
1	69.732119	6.621425	13.815199	9.825255	6.259978	11	5.198540	0.560063	10.755058	8.039328	13.6522
2	52.448930	10.773829	5.802397	30.974951	5.895796	8	4.880866	0.563916	25.677644	24.438971	12.5632
3	81.918734	1.370870	4.290594	12.413444	4.434488	7	3.417664	0.687313	2.469991	7.604859	13.0013
4	75.649980	10.455473	4.733141	9.161406	3.371602	5	2.579847	0.753247	11.680421	10.101772	12.4896

```
In [3]: data.describe()
```

```
Out[3]:
```

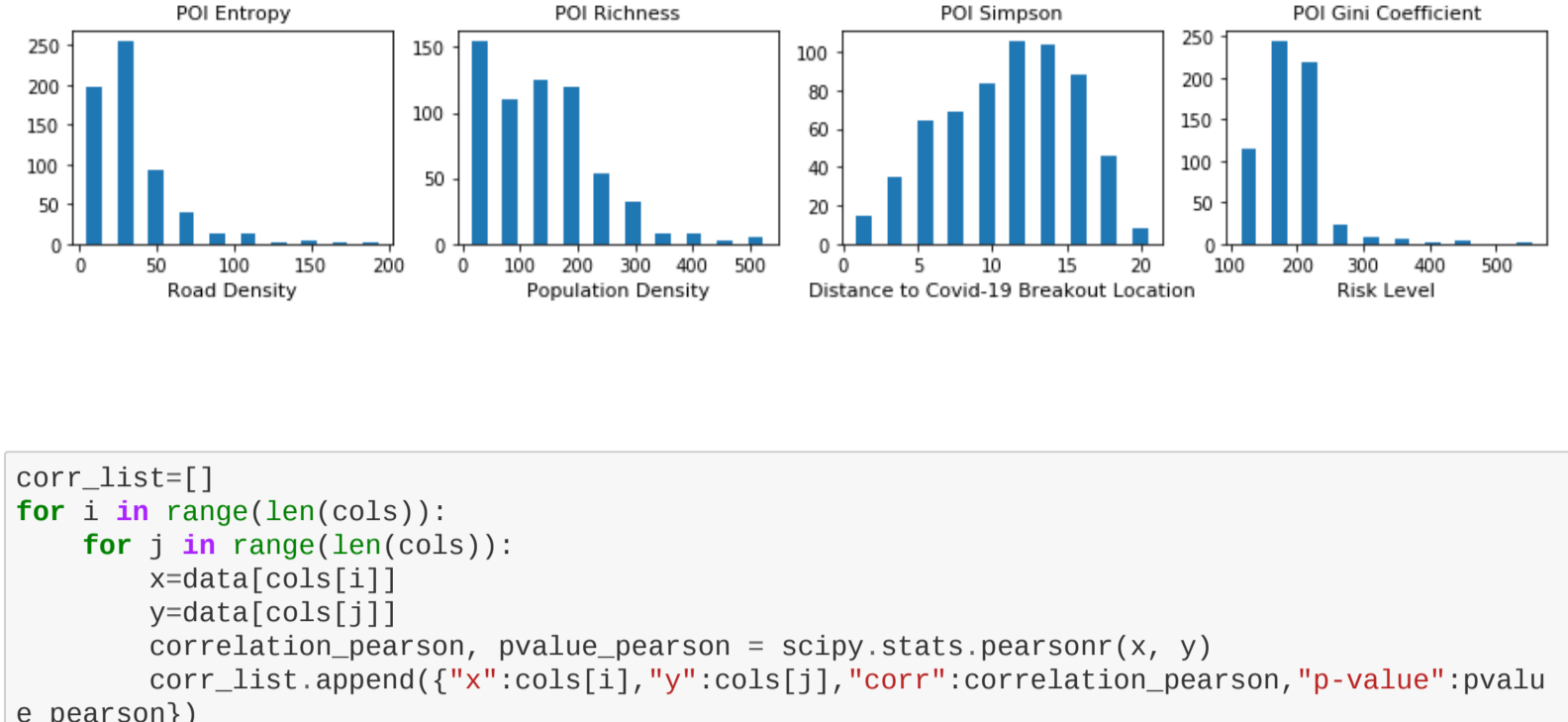
	Open Area	Traffic Area	Residential Area	Commercial Area	POI Entropy	POI Richness	POI Simpson	POI Gini Coefficient	Road Density	Population Density	Distance to Covid Breakout Location
count	619.000000	619.000000	619.000000	619.000000	619.000000	619.000000	619.000000	619.000000	619.000000	619.000000	619.000000
mean	30.343190	10.494497	42.866562	16.294614	6.706165	10.110648	6.706165	0.5458242	0.521041	33.438890	136.0594
std	26.084727	6.332472	25.502116	11.299997	1.491362	1.960616	1.522794	0.107144	25.350221	67.1282	67.1282
min	0.000000	0.000000	0.000000	0.000000	1.000000	1.000000	1.000000	0.248521	0.000000	1.5577	1.5577
25%	8.036251	6.683807	20.063942	8.064498	5.888762	10.000000	4.428854	0.449962	16.392934	55.1807	55.1807
50%	23.696392	9.509252	42.831070	14.091025	6.995776	11.000000	5.594995	0.508385	27.730754	127.3777	127.3777
75%	50.620470	12.901970	64.753794	21.873750	7.731973	11.000000	6.528810	0.583542	41.158415	167.9506	167.9506
max	100.000000	73.583656	100.000000	77.899569	9.917969	11.000000	9.198390	0.909091	198.452300	537.0502	537.0502

```
In [4]: data.columns
```

```
Out[4]: Index(['Open Area', 'Traffic Area', 'Residential Area', 'Commercial Area', 'POI Entropy', 'POI Richness', 'POI Simpson', 'POI Gini Coefficient', 'Road Density', 'Population Density', 'Distance to Covid-19 Breakout Location', 'Risk Level'], dtype='object')
```

```
In [5]: cols=data.columns.tolist()
```

```
fig,ax=plt.subplots(3,4,figsize=(15,8))
plt.subplots_adjust(wspace=0.2, hspace=0.3)
axs=ax.flatten()
for col in cols:
    _=axs[cols.index(col)].hist(data[col],width=0.5)
axs[cols.index(col)].set_xlabel(col,fontsize=11)
plt.show()
```



```
In [6]: corr_list=[]
for i in range(len(cols)):
    for j in range(len(cols)):
        x=data[cols[i]]
        y=data[cols[j]]
        correlation_pearson, pvalue_pearson = scipy.stats.pearsonr(x, y)
        corr_list.append(("x":cols[i],"y":cols[j],"corr":correlation_pearson,"p-value":pvalue_pearson))

corr_df=pd.DataFrame(corr_list)
print(" r :")
corr_df1=corr_df.pivot_table(index="x",columns="y",values="corr")
corr_df1=corr_df1.reindex(index=cols,columns=cols)
corr_df1.to_csv("pearson_corr.csv")
display(corr_df1)

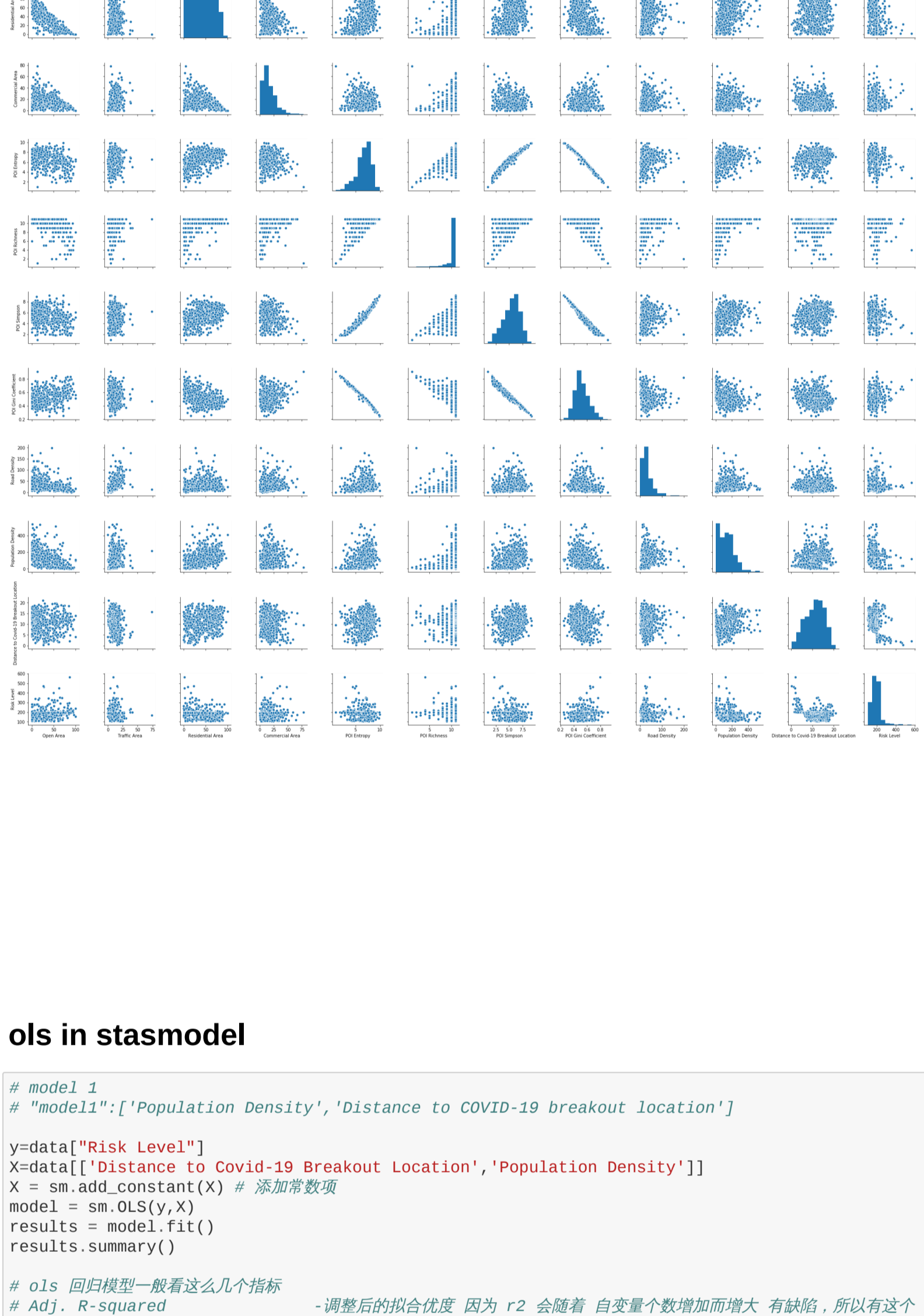
print(" p :")
corr_df2=corr_df.pivot_table(index="x",columns="y",values="p-value")
corr_df2=corr_df2.reindex(index=cols,columns=cols)
corr_df2.to_csv("pearson_corr_pvalue.csv")
display(corr_df2)
r :
```

	Open Area	Traffic Area	Residential Area	Commercial Area	POI Entropy	POI Richness	POI Simpson	POI Gini Coefficient	Road Density	Population Density
Open Area	1.000000	-0.282023	-0.844146	-0.240805	-0.390280	-0.506744	-0.328219	0.370190	-0.300703	-0.526110
Traffic Area	-0.282023	1.000000	-0.084538	0.281876	0.008710	0.097651	0.009007	-0.011634	0.501339	0.151195
Residential Area	-0.844146	-0.084538	1.000000	-0.266137	-0.286137	0.473444	0.439070	0.423015	-0.456373	0.107101
Commercial Area	-0.240805	0.281876	-0.266137	1.000000	-0.170362	0.121880	-0.204244	0.184286	0.170950	0.106199
POI Entropy	-0.390280	0.008710	0.473444	-0.170362	1.000000	0.619359	0.973512	-0.992663	0.054265	0.306204
POI Richness	-0.506744	0.097651	0.439070	0.121880	0.619359	1.000000	0.482504	-0.579733	0.033140	0.356479
POI Simpson	-0.328219	0.009007	0.423015	-0.204244	0.973512	0.482504	1.000000	-0.980355	0.049994	0.272451
POI Gini Coefficient	0.370190	-0.011634	-0.456373	0.184286	-0.992663	-0.579733	-0.980355	1.000000	-0.048731	-0.303625
Road Density	-0.300703	0.501339	0.107101	0.170950	0.054265	0.033140	0.049994	-0.048731	1.000000	0.211977
Population Density	-0.526110	0.151195	0.452489	0.106199	0.306204	0.356479	0.272451	-0.303625	0.211977	1.000000

	Open Area	Traffic Area	Residential Area	Commercial Area	POI Entropy	POI Richness	POI Simpson	POI Gini Coefficient
Open Area	0.000000e+00	8.800190e-13	3.399519e-169	1.290493e-09	5.904662e-24	1.090873e-41	5.167163e-17	1.537782e-2
Traffic Area	8.800190e-13	0.000000e+00	3.548307e-02	9.053225e-13	8.287764e-01	1.508127e-02	8.230330e-01	7.726746e-0
Residential Area	3.399519e-169	3.548307e-02	0.000000e+00	1.699390e-11	1.423510e-35	1.471276e-30	2.893681e-28	3.585719e-3
Commercial Area	1.290493e-09	9.053225e-13	1.699390e-11	0.000000e+00	2.033056e-05	2.385125e-03	2.969371e-07	3.925058e-0
POI Entropy	5.904662e-24	8.287764e-01	1.423510e-35	2.033056e-05	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
POI Richness	1.090873e-41	1.508127e-02	1.471276e-30	2.385125e-03	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
POI Simpson	5.167163e-17	8.230330e-01	2.893681e-28	2.969371e-07	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
POI Gini Coefficient	1.537782e-21	7.726746e-01	3.585719e-33	3.925058e-06	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
Road Density	2.107377e-14	1.049905e-40	7.654585e-03	1.901973e-05	1.775443e-01	4.104729e-01	2.141989e-01	2.260170e-0
Population Density	2.324844e-45	1.594847e-04	1.426754e-32	8.185143e-03	6.660735e-15	5.507127e-20	5.360644e-12	1.146320e-1
Distance to Covid-19 Breakout Location	3.222418e-02	1.988471e-07	5.764893e-07	6.846189e-04	3.371166e-05	7.603171e-02	1.487242e-05	9.809937e-0
Risk Level	6.719209e-11	9.776750e-01	1.185434e-10	6.209858e-01	1.745043e-07	5.942641e-07	1.980102e-06	2.191767e-0

```
In [7]: sns.pairplot(data)
```

```
Out[7]: <seaborn.axisgrid.PairGrid at 0x1de639783128>
```



## ols in stamodel

```
In [8]: # model 1
# "model1":["Population Density", 'Distance to COVID-19 breakout location']
y=data["Risk Level"]
X=data[["Distance to Covid-19 Breakout Location", 'Population Density']]
X = sm.add_constant(X) # 添加常数项
model = sm.OLS(y,X)
results = model.fit()
results.summary()
```

```
Out[8]:
```

Dep. Variable:	Risk Level	R-squared:	0.167				
Model:	OLS	Adj. R-squared:	0.165				
Method:	Least Squares	F-statistic:	61.84				
Date:	Thu, 27 Aug 2020	Prob (F-statistic):	3.35e-25				
Time:	18:15:57	Log-Likelihood:	-3212.7				
No. Observations:	619	AIC:	6431.				
Df Residuals:	616	BIC:	6445.				
Df Model:	2						
Covariance Type:	nonrobust						
	coef	std err	t	P> t	[0.025	0.975]	
	const	239.6807	5.047	47.491	0.000	229.769	249.592
Distance to Covid-19 Breakout Location	-3.4914	0.420	-8.316	0.000	-4.316	-2.667	
Population Density	-0.0962	0.019	-5.180	0.000	-0.133	-0.060	
Omnibus:	332.587	Durbin-Watson:	1.750				
Prob(Omnibus):	0.000	Jarque-Bera (JB):	3561.330				
Skew:	2.156	Prob(JB):	0.00				
Kurtosis:	13.931	Cond. No.	484.				

```
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
```

```
In [9]: # model 2
# "model1":["Population Density", 'Distance to COVID-19 breakout location', 'Road density']
y=data["Risk Level"]
X=data[["Distance to Covid-19 Breakout Location", 'Population Density', 'Road Density']]
X = sm.add_constant(X) # 添加常数项
model = sm.OLS(y,X)
results = model.fit()
results.summary()
```

```
Out[9]:
```

Dep. Variable:	Risk Level	R-squared:	0.167				
Model:	OLS	Adj. R-squared:	0.163				
Method:	Least Squares	F-statistic:	41.22				
Date:	Thu, 27 Aug 2020	Prob (F-statistic):	2.80e-24				
Time:	18:19:16	Log-Likelihood:	-3212.6				
No. Observations:	619	AIC:	6433.				
Df Residuals:	615	BIC:	6451.				
Df Model:	3						
Covariance Type:	nonrobust						
	coef	std err	t	P> t	[0.025	0.975]	
	const	238.9630	5.047	44.196	0.000	228.345	249.581
Distance to Covid-19 Breakout Location	-3.4875	0.420	-8.298	0.000	-4.313	-2.662	
Population Density	-0.0977	0.019	-5.138	0.000	-0.135	-0.060	
Road Density	0.0263	0.070	0.372	0.710	-0.113	0.165	
Omnibus:	332.230	Durbin-Watson:	1.750				
Prob(Omnibus):	0.000	Jarque-Bera (JB):	3552.148				
Skew:	2.153	Prob(JB):	0.00				
Kurtosis:	13.917	Cond. No.	526.				

```
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
```

```
In [10]: # model 3
# "model1":["Population Density", 'Distance to COVID-19 breakout location', 'POI Entropy', 'POI Richness', 'POI Simpson', 'Open Area', 'Traffic Area', 'Residential Area', 'Commercial Area']
y=data["Risk Level"]
X=data[["Distance to Covid-19 Breakout Location", 'Population Density', 'POI Entropy', 'POI Richness', 'POI Simpson', 'Open Area', 'Traffic Area', 'Residential Area', 'Commercial Area']]
X = sm.add_constant(X) # 添加常数项
model = sm.OLS(y,X)
results = model.fit()
results.summary()
```

```
Out[10]:
```

Dep. Variable:	Risk Level	R-squared:	0.182				
Model:	OLS	Adj. R-squared:	0.176				
Method:	Least Squares	F-statistic:	27.32				
Date:	Thu, 27 Aug 2020	Prob (F-statistic):	5.38e-25				
Time:	18:22:35	Log-Likelihood:	-3207.1				
No. Observations:	619	AIC:	6426.				
Df Residuals:	613	BIC:	6453.				
Df Model:	5						
Covariance Type:	nonrobust						
	coef	std err	t	P> t	[0.025	0.975]	
	const	275.7921	12.321	22.383	0.000	251.959	299.989
Distance to Covid-19 Breakout Location	-3.4484	0.421	-8.191	0.000	-4.275	-2.622	
Population Density	-0.0729	0.020	-3.684	0.000	-0.112	-0.034	
POI Entropy	-4.0281	7.679	-0.525	0.600	-19.003	11.052	
POI Richness	-2.3464	1.955	-1.200	0.231	-6.187	1.494	
POI Simpson	2.0149	6.748	0.299	0.765	-11.237	15.267	
Omnibus:	311.290	Durbin-Watson:	1.752				
Prob(Omnibus):	0.000	Jarque-Bera (JB):	2854.821				
Skew:	2.033	Prob(JB):	0.00				
Kurtosis:	12.703	Cond. No.	1.20e+03				

```
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.2e+03. This might indicate that there are strong multicollinearity or other numerical problems.
```

```
In [11]: # model 4
# "model1":["Population Density", 'Distance to COVID-19 breakout location', 'POI Entropy', 'POI Richness', 'POI Simpson', 'Open Area', 'Traffic Area', 'Residential Area', 'Commercial Area', 'Road Density']
y=data["Risk Level"]
X=data[["Distance to Covid-19 Breakout Location", 'Population Density', 'POI Entropy', 'POI Richness', 'POI Simpson', 'Open Area', 'Traffic Area', 'Residential Area', 'Commercial Area', 'Road Density']]
X = sm.add_constant(X) # 添加常数项
model = sm.OLS(y,X)
results = model.fit()
results.summary()
```

```
Out[11]:
```

Dep. Variable:	Risk Level	R-squared:	0.182			
Model:	OLS	Adj. R-squared:	0.174			
Method:	Least Squares	F-statistic:	22.73			
Date:	Thu, 27 Aug 2020	Prob (F-statistic):	2.96e-24			
Time:	18:24:27	Log-Likelihood:	-3209.1			
No. Observations:	619	AIC:	6428.			
Df Residuals:	609	BIC:	6459.			
Df Model:	6					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
	const	284.5578	153.565			