

Project 7. Geographically Weighted Regression (Spatial Modeling) of Covid 19 Cases in the City of Chicago

Abstract

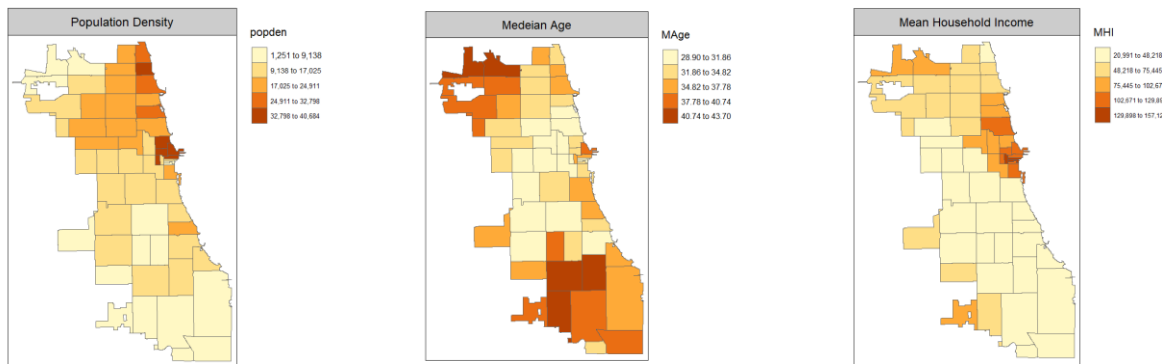
This study used Geographically Weighted Regression (GWR) for spatial analysis of Covid-19 cases in Chicago, integrating weekly test data with demographic variables. The study analyzes how test frequency was affected by population density, median age, and mean household income across different zip codes. OLS and two types of GWR packages were used for the analysis. The result showed significantly high residuals, a possible reason that can be related to the variance of the total test done in each area.

1. Data Exploration and Selection of Interest Factors

The main task of this section is to explore the content of weekly “cases” and demographic data. For this project, I selected Weekly Test (shown in “Weekly...Test”) as my interest dependent variable (y) as it represents the test done every week. For the independent variables, I selected Population Density (shown in “popden”), Mean House Income (shown in “MHI”), and Median Age (shown in “MAge”) in the demographic data. The first step of this section is to merge the data. There are three parts of data: Chicago boundaries based on zip code, Chicago Covid-19 cases, and Chicago demographics. All three files will be merged based on the zip codes (shown in “zip”). The Figure 01 shows the thematic map of three attributes.

Figure 01

Thematic Map of Three Attributes



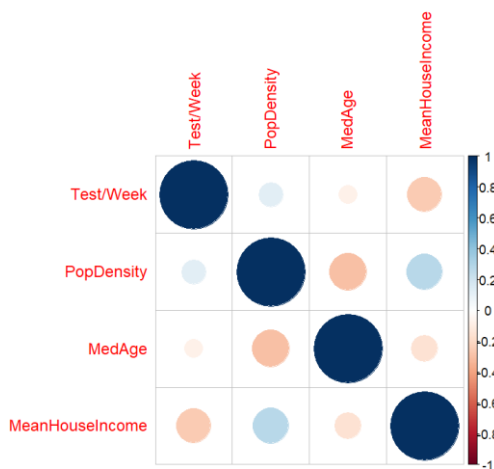
The first map shows the distribution of population density across different areas. The color gradient indicates that the population density is highest in the northeastern part, with some areas having a density as high as 32,798 to 40,684 people per zip code. The density in southern Chicago is relatively lower than northern part of Chicago. The second map illustrates the median age of residents in each area. The darker regions, particularly in the southeast and northwest, indicate a higher median age, reaching between 40.74

to 43.70 years. Conversely, some northeastern and central areas have a younger demographic, with median ages as low as 28.90 to 31.86 years. The third map depicts the average household income. There's a noticeable concentration of higher incomes (102,671 to 157,125 USD) in the Chicago downtown area. In contrast, the southern areas have lower mean household incomes, ranging from 20,991 to 48,218 USD.

After the exploration, a correlation of independent variables and animation of weekly cases were conducted to find out the strength of correlations between variables and the change of weekly cases. Figure 02 shows the correlation matrix of variables.

Figure 02

Correlation Matrix



For population density and median age, there is a slight negative correlation between these two variables (-0.4), as indicated by the light red circle in the matrix. This suggests that areas with higher population densities tend to have a younger population, whereas areas with lower densities tend to have an older median age. For population density and Mean House Income, the correlation between these two variables appears to be slightly positive (blue +0.4). This implies that areas with higher population densities might have higher mean household incomes. For median age and Mean House Income, there seems to be a weak negative correlation between median age and mean household income (light

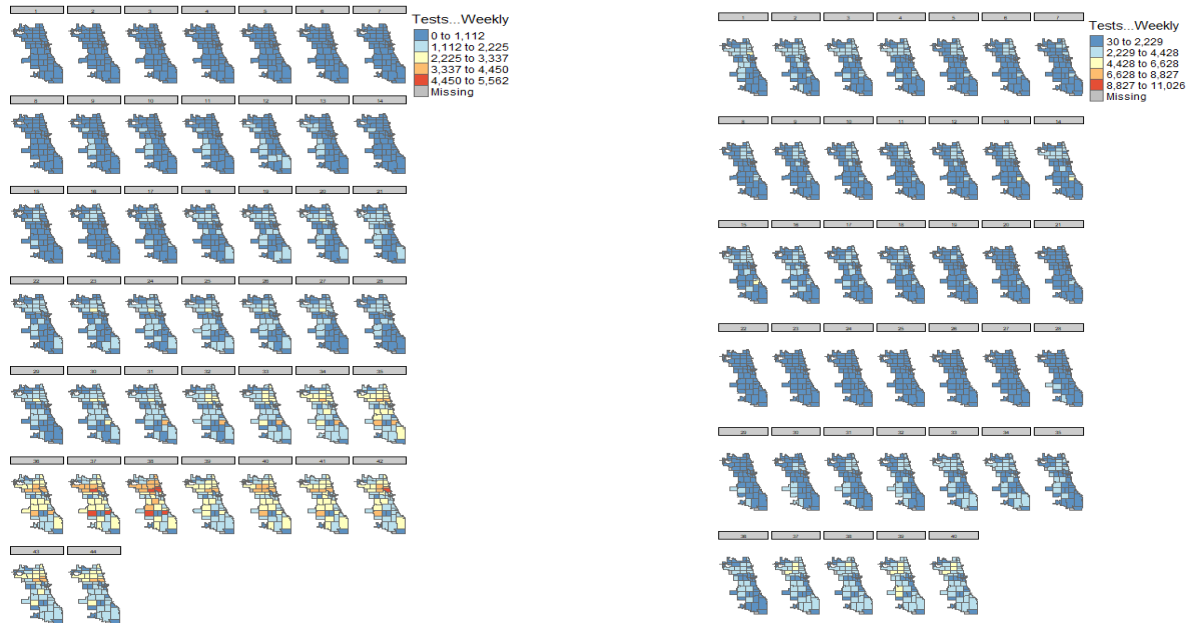
orange -0.2), indicating that areas with an older median age might have slightly lower mean household incomes.

For the animation map of cases, since the data ranges from first week of March in 2020 to End of October in 2021, the cases data was split into two years first before the animation map. The case data in 2020 contains 44 weeks (March 01 to end of 2020), and the data in 2021 contains 40 weeks (January 03 to October 09, 2021). Figure 03 shows the animation map of the weekly test in 2020 (left) and 2021(right). It is clear to see that since the spread of Covid-19 in the United States, the weekly test significantly increased in the mid-2020. However, it is interesting to see that starting from 22nd week of 2021 (May), the weekly test decreased during summer time, one possible reason can be since the temperature getting warmer, less people would get Covid-19 compared to winter time.

Another possible reason can be the vaccine deployment as most people are eligible to get Covid-19 vaccine during that time, fewer people would get infected, resulting in fewer Covid-19 weekly tests.

Figure 03

Animation Map of Weekly Test in 2020 and 2021.

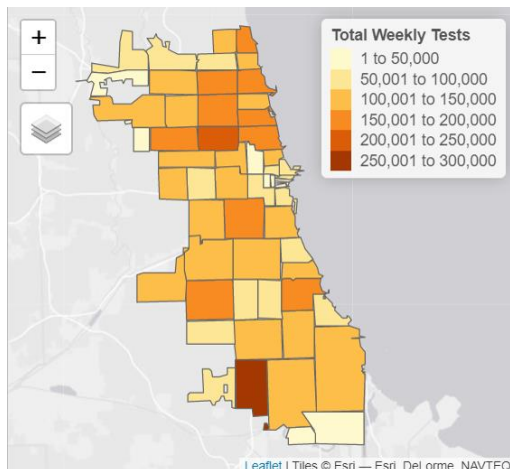


2. Ordinary Linear Model and Spatial Distribution

The first step of this section is to get the aggregated weekly test data. After summing up the weekly test data based on each zip code, the thematic map (Figure 04) shows the distribution of total test conducted in each zip code.

Figure 04

Total Test by Each Zip Code



It is surprising to see the variation in the total Covid-19 tests done by each zip code area, as some of the county have conducted less than 50,000 test while others have done more than 250,000 tests in total. The significant differences in the total test between each zip code can cause significant high or low residuals in the following test. One possible reason can be the testing facilities are located at certain areas in Chicago, which caused most of the people in the city of Chicago to conduct the Covid-19 test at certain locations, while other areas may have fewer testing sites around.

For the ordinary linear model (OLS), the regression formula is shown below:

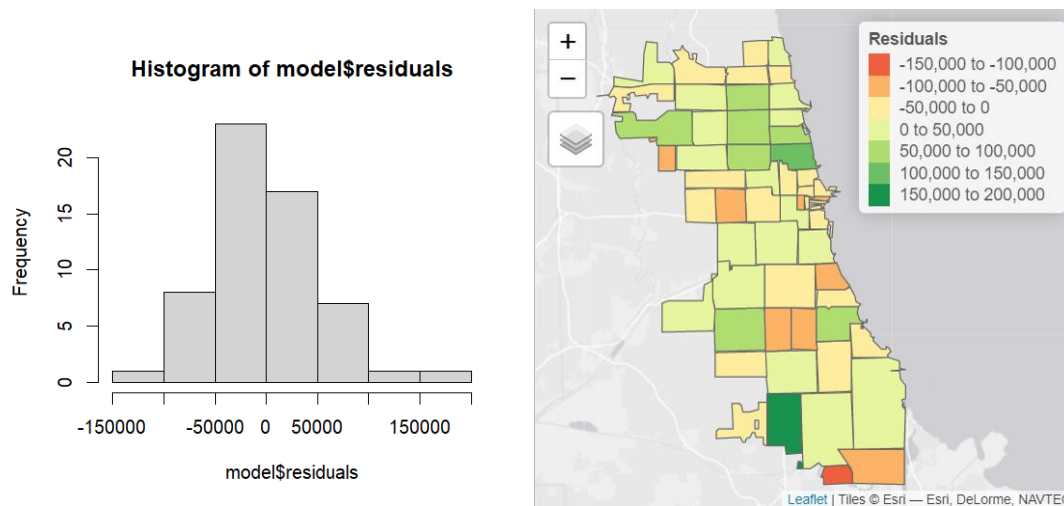
Total_Test_Weekly ~ PopDensity + MedAge + MeanHouseIncome.

Data = aggregated_data.

The result showed that only two variables showed statistically significant results: the population density showed one star ($P < 0.05$) next to it, and the MHI showed two stars ($P < 0.01$) next to it. The R-squared showed that the model can explain 19.8% of the data, indicating that other factors were not selected, which can affect the result. The screenshot of the summary is shown in the reference. From the residual plot, the result showed that observations 14, 41, 58 do not follow the trend, which can be the outliers of the data. The variance of the residual is also not constant, and some of the points indicated a high Cook's distance (reference). Figure 05 shows the spatial distribution and histogram of residuals.

Figure 05

Histogram and Spatial Distribution of Residuals



It is clear to see that there are some outliers with significant high and low residuals in the map. At the southern end of the map, it is clear to see the difference in the residuals (green and red), which proved from the total test (Figure 04 above) that the significant difference in the test caused the huge residuals.

3. GWR (spgwr Package)

The “spgwr” package was used to conduct the geographical weight regression analysis in this section. The initial process is to find the proper bandwidth. The GWR with three different kernel functions (geographical weight) were used for this section. Figure 06 shows the equation of geographical weight.

Figure 06

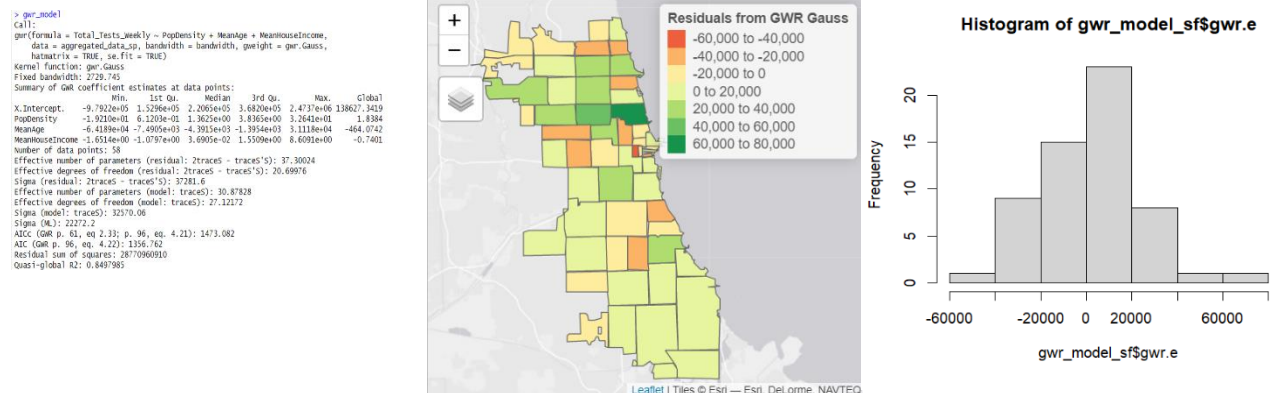
Functions of GWR

Gaussian	$w_{ij} = \exp\left(-\frac{1}{2}\left(\frac{d_{ij}}{b}\right)^2\right)$
Bisquare	$w_{ij} = \begin{cases} (1 - (d_{ij}/b)^2)^2 & \text{if } d_{ij} < b \\ 0 & \text{otherwise} \end{cases}$
Tricube	$w_{ij} = \begin{cases} (1 - (d_{ij} /b)^3)^3 & \text{if } d_{ij} < b \\ 0 & \text{otherwise} \end{cases}$

The fixed bandwidth for the Gaussian Kernel model is 2,729.745 meters, indicating the scale of spatial variation for which the model accounts. The model selection criteria are indicated by an AIC (Akaike Information Criterion) of 1356.762 and an AICc (corrected Akaike Information Criterion for small sample sizes) of 1473.082. The Quasi-global R^2 value is 0.8497985, suggesting that approximately 84.98% of the variation in the dependent variable is explained by the model across the study area, which is a robust level of explanatory power. Figure 07 shows the result of the GWR using Gaussian Kernel, histogram, and spatial distribution of the residuals. It is clear that with the weight variable's implementation, the residual's size significantly decreased.

Figure 07

Results of GWR (Gauss)



The fixed bandwidth for the Bisquare Kernel model is 25,545.67 meters, indicating the scale of spatial variation for which the model accounts. The model selection criteria are indicated by an AIC of 1429.68 and an AICc (corrected Akaike Information Criterion for small sample sizes) of 1417.68. The Quasi-global R^2 value is 0.354333, suggesting that approximately 35.43% of the variation in the dependent variable is explained by the model across the study area; compared to the Gaussian Kernel, there is a significant drop in the percentage. In addition, there is a significant increase of the fixed bandwidth (results in appendix).

The fixed bandwidth for the Tricube Kernel model is 25,640.3 meters, indicating the scale of spatial variation that the model accounts for. The model selection criteria are indicated by an AIC of 1430.713 and an AICc (corrected Akaike Information Criterion for small sample sizes) of 1419.464. The Quasi-global R^2 value is 0.3270923, suggesting that approximately 32.71% of the variation in the dependent variable is explained by the model across the study area; compared to the Gaussian Kernel, there is a significant drop of the percentage. Compared to Bisquare Kernel, the results are close, one possible reason is that the formula (Figure 06) for Bisquare and Tricube are similar (results in appendix).

4. GWR in Different Package

A different package, “GWModel,” was used for weight regression analysis. similar to the “spgwr” package above, three different models were used for the analysis. For the Gaussian Kernel mode, the adaptive bandwidth for this GWR model is 58. The AIC value is 1377.205, and the AICc value is 1448.112. The R-square for the model indicates that 76.25% of the data can be explained by the model, which is considered good of fit, but it is less than the result from “spgwr” package. For the Bisquare model, the result showed that the adaptive bandwidth is 58. The AIC value is 1417.86, the AICc value is 1429.945. The R-square for the model indicates that the model can explain 35.08% of the data, and the results are lower than the results in the Gaussian Kernel Model. For the Tricube model, the adaptive bandwidth is 58, The AIC value is 1418.339, and the AICc value is 1430.078. The R-square for the model indicates that the model can explain 34.32% of the data. The results are close to the bisquare results but are lower than the Gaussian Kernel Model. Table 01 shows the difference between the results from the two packages.

Table 01*Spgwr vs. GWModel Results*

Method	Bandwidth	AIC	AICc	R-square
Spgwr.Gauss	2,729.745	1356.762	1473.082	84.98%
Spgwr.Bisquare	25,545.67	1429.68	1417.68	35.43%
Spgwr.Tricube	25,640.3	1430.713	1419.464	32.71%
GWM.Gauss	58	1377.205	1448.112	76.25%
GWM.Bisquare	58	1417.86	1429.945	35.08%
GWM.Tricube	58	1418.339	1430.078	34.32%

For AIC and AICc, across all kernel methods, the spgwr package tends to produce lower AIC and AICc values compared to the GWmodel package. Lower AIC and AICc values generally suggest a better model fit, in this project, it seems that the spgwr models fit the data better. However, it can be different for different selected attributes and dependent variables. The R-squared values, which indicate the proportion of variance explained by the model, are consistently higher in the spgwr package results. For example, the Gaussian kernel has an R-squared of 84.98% in spgwr compared to 76.25% in GWmodel. This suggests that the spgwr package models, at least for those attribute and dependent variable selections and with these settings, explain more variance than those from the GWmodel package. The Bisquare and Tricube kernels show similar results to each other but are less effective than the Gaussian kernel, according to these results.

5. Effect of Different Bandwidth

The spgwr package shows a significant difference in bandwidths for Gaussian versus Bisquare and Tricube kernels. A smaller bandwidth for the Gaussian kernel might indicate that it is capturing more local variation in the data, which is reflected in the higher R-squared value. This suggests that the Gaussian kernel provides a more nuanced model of local effects. In contrast, larger bandwidths for the Bisquare and Tricube kernels could mean that these models cover a larger area, possibly ignoring local effects, reflected in the lower R-squared values. The GWModel may also use cross-validation or an information criterion to select optimal bandwidth. However, it does not use the distance, rather, it uses the number of neighbors to indicate the bandwidth. This is why the bandwidth in GWModel are all the same. However, in areas with outlier data, especially at the edges of the selected data, the

method of using neighbor as bandwidth may potentially skew the results. Returning to the thematic map of the Total Weekly Test, we can see that some edge zip codes (at the southern edge) have significantly high/low numbers of tests. Those can significantly affect the result, which is why the GWModel produces lower results than the spgwr model.

6. Summary/Conclusion/Concluding Remarks

The study conducted OLS, and GWR (two packages) analysis of the weekly test (dependent variable) along with population density, median age, and mean house income (independent variables). The result showed significant high residuals for the analysis. One possible reason may be due to the variance of the test cases in different zip codes. While some zip codes have nearly 200,000 tests, others may have less than 800 tests. Different gwr packages may have different pros and cons while doing the analysis; in addition, the selection of the bandwidths can also affect the result of whether the data fits.

Acknowledgement

I would like to thank my classmates and professor for pointing out the huge residual issue. For the secondary analysis, it turns out that it is due to the variance total test between different zip codes. To have a better explanation, a thematic map was added before the OLS test for better clarification.

References

OLS

```
> summary(model)
```

Call:
lm(formula = Total_Tests_Weekly ~ PopDensity + MeanAge + MeanHouseIncome,
data = aggregated_data)

Residuals:

	Min	1Q	Median	3Q	Max
	-106486	-34466	-1734	29183	172254

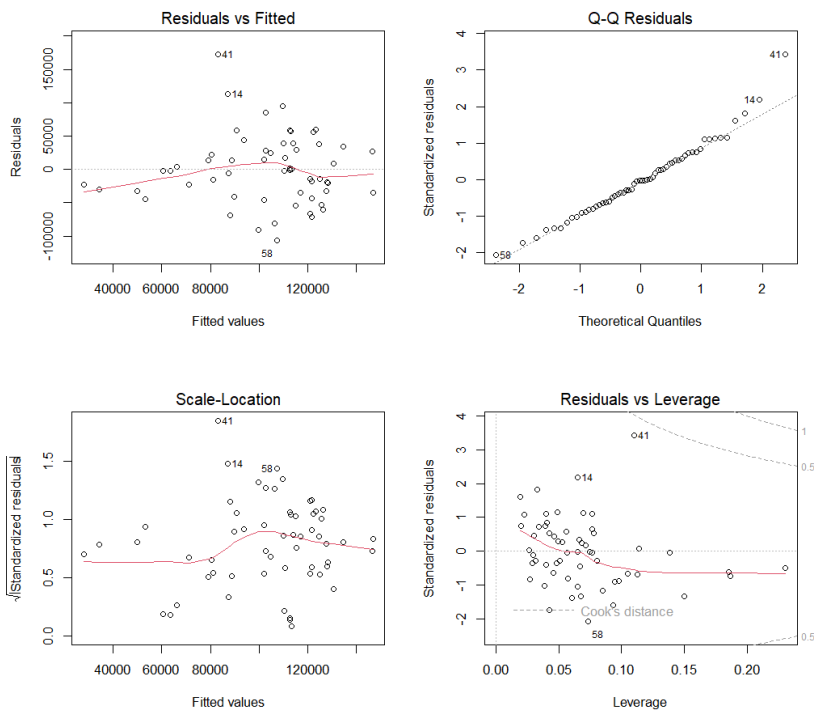
Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.386e+05	7.405e+04	1.872	0.06661 .
PopDensity	1.838e+00	8.081e-01	2.275	0.02691 *
MeanAge	-4.641e+02	1.933e+03	-0.240	0.81114
MeanHouseIncome	-7.401e-01	2.223e-01	-3.329	0.00157 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 53340 on 54 degrees of freedom
Multiple R-squared: 0.198, Adjusted R-squared: 0.1534
F-statistic: 4.443 on 3 and 54 DF, p-value: 0.007315

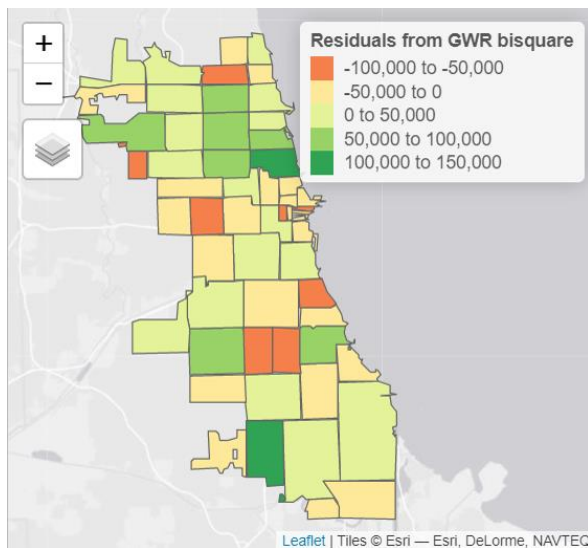
```
> # Extract residuals
```



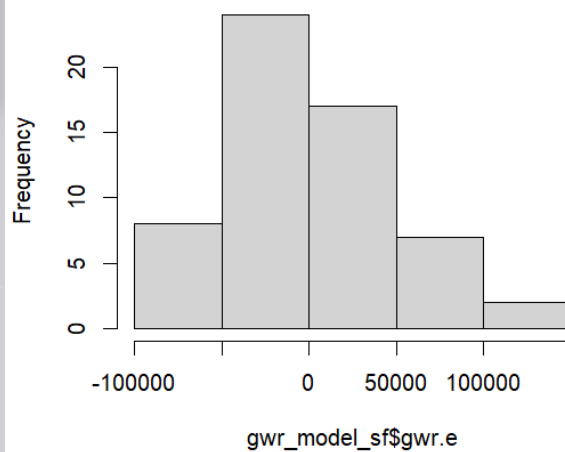
GWR

Bisquare

```
> gwr_model
Call:
gwr(formula = Total_Tests_Weekly ~ PopDensity + MeanAge + MeanHouseIncome,
     data = aggregated_data_sf, bandwidth = bandwidth, gweight = gwr.bisquare,
     hatmatrix = TRUE, se.fit = TRUE)
Kernel function: gwr.bisquare
Fixed bandwidth: 25545.67
Summary of GWR coefficient estimates at data points:
              Min.      1st Qu.      Median      3rd Qu.      Max.      Global
X.Intercept. -2.3201e+05  1.4286e+05  2.2390e+05  2.3467e+05  2.5058e+05 138627.3419
PopDensity    1.0910e+00  1.5054e+00  1.6712e+00  1.7666e+00  7.3840e+00   1.8384
MeanAge       -3.1965e+03 -2.8876e+03 -2.6391e+03 -3.3247e+02  6.9816e+03 -464.0742
MeanHouseIncome -9.1800e-01 -8.9260e-01 -8.2943e-01 -7.3054e-01  3.5731e-01  -0.7401
Number of data points: 58
Effective number of parameters (residual: 2traces - traces'S): 8.567169
Effective degrees of freedom (residual: 2traces - traces'S): 49.43283
Sigma (residual: 2traces - traces'S): 50019.19
Effective number of parameters (model: traces): 7.214035
Effective degrees of freedom (model: traces): 50.78596
Sigma (model: traces): 49348.33
Sigma (ML): 46177.49
AICc (GWR p. 61, eq 2.33; p. 96, eq. 4.21): 1429.996
AIC (GWR p. 96, eq. 4.22): 1417.68
Residual sum of squares: 123676935628
Quasi-global R2: 0.354333
```



Histogram of gwr_model_sf\$gwr.e



Tricube

Call:

```
gwr(formula = Total_Tests_Weekly ~ PopDensity + MeanAge + MeanHouseIncome,
     data = aggregated_data_sf, bandwidth = bandwidth, gweight = gwr.tricube,
     hatmatrix = TRUE, se.fit = TRUE)
```

Kernel function: gwr.tricube

Fixed bandwidth: 26408.3

Summary of GWR coefficient estimates at data points:

	Min.	1st Qu.	Median	3rd Qu.	Max.	Global
X.Intercept.	-2.0001e+05	1.3711e+05	2.2258e+05	2.3176e+05	2.4649e+05	138627.3419
PopDensity	1.1179e+00	1.5097e+00	1.6823e+00	1.7413e+00	6.9471e+00	1.8384
MeanAge	-3.1358e+03	-2.8217e+03	-2.6283e+03	-1.6162e+02	6.2618e+03	-464.0742
MeanHouseIncome	-9.1221e-01	-8.7012e-01	-8.3836e-01	-7.4161e-01	3.1159e-01	-0.7401

Number of data points: 58

Effective number of parameters (residual: 2traces - traces'S): 7.504417

Effective degrees of freedom (residual: 2traces - traces'S): 50.49558

Sigma (residual: 2traces - traces'S): 50523.23

Effective number of parameters (model: traces): 6.601686

Effective degrees of freedom (model: traces): 51.39831

Sigma (model: traces): 50077.59

Sigma (ML): 47141.55

AICc (GWR p. 61, eq. 2.33; p. 96, eq. 4.21): 1430.713

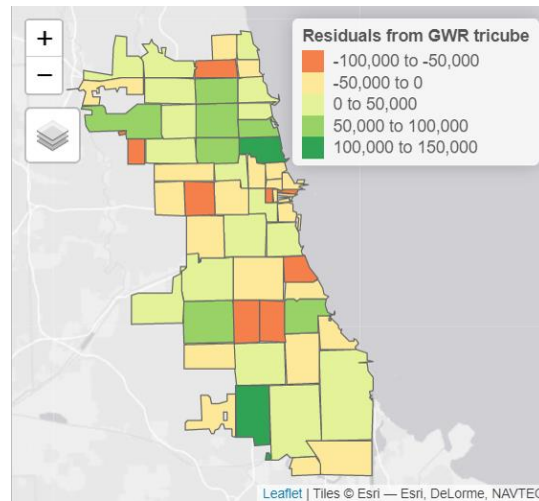
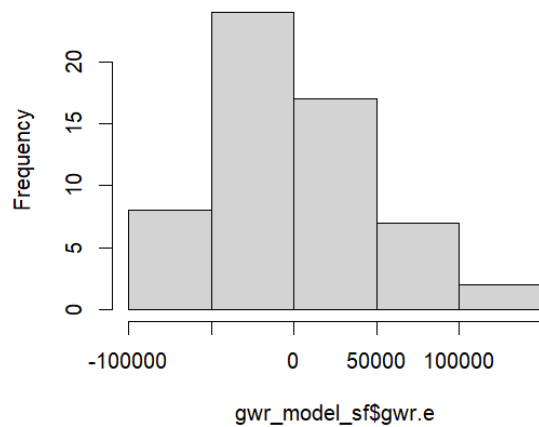
AIC (GWR p. 96, eq. 4.22): 1419.464

Residual sum of squares: 128894875378

Quasi-global R2: 0.3270923

...

Histogram of gwr_model_sf\$gwr.e



GW Model Gaussian

```

-----
*               Results of Global Regression               *
-----
Call:
lm(formula = formula, data = data)

Residuals:
    Min       1Q   Median       3Q      Max
-106486  -34466  -1734    29183   172254

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.386e+05  7.405e+04   1.872  0.06661 .
PopDensity   1.838e+00  8.081e-01   2.275  0.02691 *
MedAge      -4.641e+02  1.933e+03  -0.240  0.81114
MeanHouseIncome -7.401e-01  2.223e-01  -3.329  0.00157 **

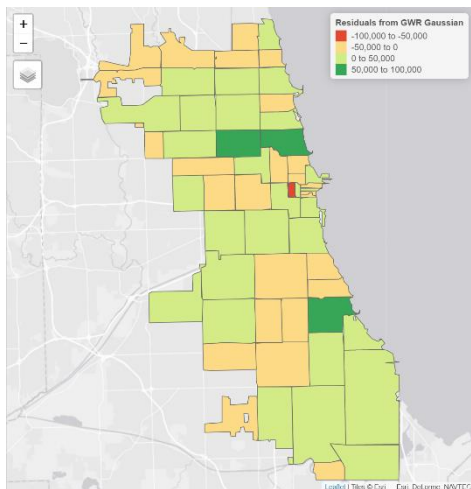
---Significance stars---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 53340 on 54 degrees of freedom
Multiple R-squared:  0.198
Adjusted R-squared:  0.1534
F-statistic: 4.443 on 3 and 54 DF,  p-value: 0.007315
***Extra Diagnostic information***
Residual sum of squares: 1.53629e+11
Sigma(hat): 52377.24
AIC: 1433.044
AICC: 1434.198
BIC: 1405.648
-----
*               Results of Geographically weighted Regression               *
-----

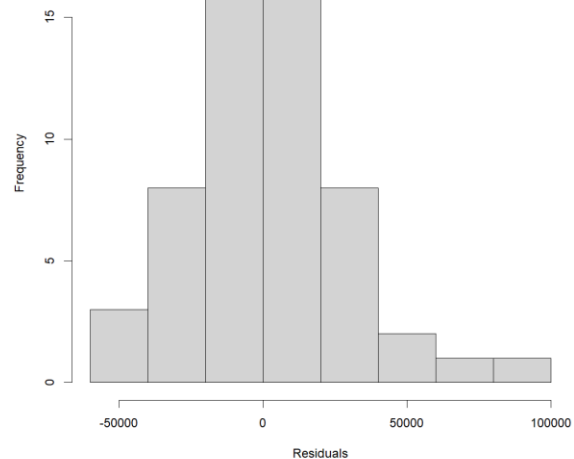
*****Model calibration information*****
Kernel function: bisquare
Fixed bandwidth: 8754.382
Regression points: the same locations as observations are used.
Distance metric: Euclidean distance metric is used.

*****Summary of GWR coefficient estimates:*****
              Min.      1st Qu.      Median      3rd Qu.      Max.
Intercept   -8.8162e+05  1.5673e+05  2.5032e+05  3.4504e+05  2.4020e+06
PopDensity   -1.6509e+01  6.3149e-01  1.2266e+00  3.1825e+00  3.7506e+01
MedAge       -6.2490e+04  -7.4588e+03  -4.2847e+03  -2.4722e+03  2.8973e+04
MeanHouseIncome -1.7309e+00 -1.0746e+00  1.1049e-01  1.3836e+00  1.2581e+01
*****Diagnostic information*****
Number of data points: 58
Effective number of parameters (2*trace(S) - trace(S'S)): 29.36872
Effective degrees of freedom (n-2*trace(S) + trace(S'S)): 28.63128
AICC (GWR book, Fotheringham, et al. 2002, p. 61, eq 2.33): 1448.112
AIC (GWR book, Fotheringham, et al. 2002, GWR p. 96, eq. 4.22): 1377.205
BIC (GWR book, Fotheringham, et al. 2002, GWR p. 61, eq. 2.34): 1394.99
Residual sum of squares: 45479741910
R-square value: 0.7625688
Adjusted R-square value: 0.5102079
-----

```



Histogram of GWR Residuals (Gaussian)



```

*****
Results of Global Regression
*****

Call:
lm(formula = formula, data = data)

Residuals:
    Min       1Q   Median       3Q      Max
-106486   -34466   -1734   29183  172254

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.388e+05  7.405e+04  1.872  0.06661
PopDensity   1.838e+00  8.081e-01  2.275  0.02691 *
MedAge      -4.641e-02  1.933e-03  -0.240  0.81114
MeanHouseIncome -7.401e-01  2.223e-01  -3.329  0.00157 **

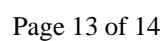
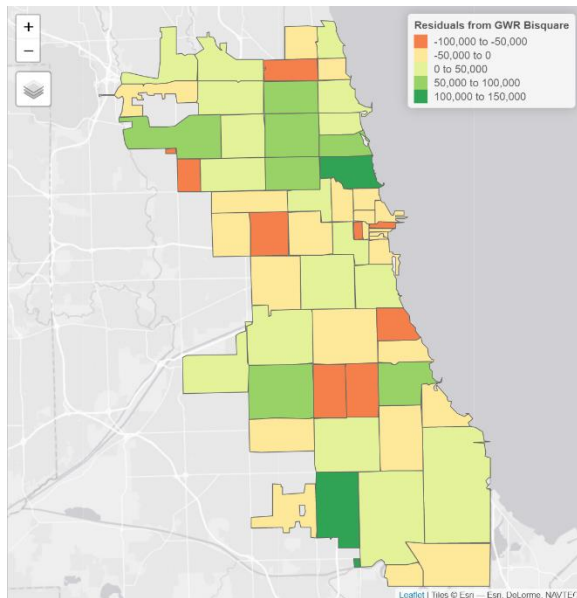
---Significance stars
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 53340 on 54 degrees of freedom
Multiple R-squared: 0.198
Adjusted R-squared: 0.1534
F-statistic: 4.443 on 3 and 54 DF,  p-value: 0.007315
***extra diagnostic information
Residual sum of squares: 1.53629e+11
Sigma(hat): 52377.24
AIC: 1433.044
AICC: 1434.198
BIC: 1405.648
*****
*               Results of Geographically weighted Regression               *
*****

*****Model calibration information*****
kernel function: bisquare
Fixed bandwidth: 26120.46
Regression points: the same locations as observations are used.
Distance metric: Euclidean distance metric is used.

*****Summary of GWR coefficient estimates:*****
              Min.           1st Qu.           Median           3rd Qu.           Max.
Intercept    -2.2224e+05  1.4213e+05  2.2085e+05  2.3239e+05  2.4817e+05
PopDensity    1.1232e+00  1.5053e+00  1.6724e+00  1.7494e+00  1.7898e+00
MedAge       -3.0891e-03  -2.8288e-03  -2.6469e-03  -3.2173e-02  6.8124e-03
MeanHouseIncome -9.1194e-01  -8.6790e-01  -8.2783e-01  -7.3217e-01  -2.6680e-01
*****extra diagnostic information*****
Number of data points: 58
Effective number of parameters (2*trace(S) - trace(S'S)): 8.39162
Effective degrees of freedom (n-2*trace(S) + trace(S'S)): 49.60838
AIC (GWR book, Fotheringham, et al. 2002, p. 61, eq 2.33): 1429.945
AICc (GWR book, Fotheringham, et al. 2002, GWR p. 96, eq. 4.22): 1417.86
BIC (GWR book, Fotheringham, et al. 2002, GWR p. 61, eq. 2.34): 1381.537
Residual sum of squares: 124343858020
R-square value: 0.3508513
Adjusted R-square value: 0.238784
*****

```



Tricube

```

*****
Results of Global Regression
*****

Call:
lm(formula = formula, data = data)

Residuals:
    Min       1Q   Median       3Q      Max
-106486  -34466  -1734   29183  172254

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.386e+05  7.405e+04   1.872  0.06661 .
PopDensity   1.838e+00  8.081e-01   2.275  0.02691 *
MedAge      -4.641e+02  1.933e+03  -0.240  0.81114
MeanHouseIncome -7.401e-01  2.223e-01  -3.329  0.00157 **

---Significance stars
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 53340 on 54 degrees of freedom
Multiple R-squared:  0.198
Adjusted R-squared:  0.1534
F-statistic: 4.443 on 3 and 54 DF, p-value: 0.007315
***Extra Diagnostic information
Residual sum of squares: 1.53629e+11
Sigma(hat): 52377.24
AIC: 1433.044
AICC: 1434.198
BIC: 1405.648
*****

Results of Geographically Weighted Regression
*****

*****Model calibration information*****
Kernel function: bisquare
Fixed bandwidth: 27000.411
Regression points: the same locations as observations are used.
Distance metric: Euclidean distance metric is used.

*****Summary of GWR coefficient estimates:*****
            Min.      1st Qu.      Median      3rd Qu.      Max.
Intercept  -2.0579e+05  1.4126e+05  2.1850e+05  2.2985e+05  2.4494e+05
PopDensity   1.1943e+00  1.5086e+00  1.6455e+00  1.7402e+00  6.5277e+00
MedAge      -3.0124e+03  -2.7503e+03  -2.5557e+03  -3.1129e+02  6.7019e+03
MeanHouseIncome -9.0333e-01 -8.8122e-01 -8.2536e-01 -7.3554e-01  1.0320e-01
*****Diagnostic information*****
Number of data points: 58
Effective number of parameters (2*trace(S) - trace(S'S)): 8.138481
Effective degrees of freedom (n-2*trace(S) + trace(S'S)): 49.86152
AICC (GWR book, Fotheringham, et al. 2002, p. 61, eq 2.33): 1430.078
AIC (GWR book, Fotheringham, et al. 2002, GWR p. 96, eq. 4.22): 1418.339
BIC (GWR book, Fotheringham, et al. 2002, GWR p. 61, eq. 2.34): 1381.412
Residual sum of squares: 125799847376
R-square value: 0.3432502
Adjusted R-square value: 0.2338605
*****

```

