

# *The Association between Urbanization and Northern Cardinal Observation in United States Counties*

By Ian Brandenburg – 2304791 (GitHub: [https://github.com/landrewburg/DA2\\_TERM\\_PROJECT](https://github.com/landrewburg/DA2_TERM_PROJECT))

## **1. Introduction**

The Northern Cardinal (*Cardinalis cardinalis*) is one of the most observed bird species in the U.S., commonly found in backyards and well-recognized nationwide<sup>1</sup>. It is referred to as a ‘backyard’ bird species, which means it is very commonly found in rural and urbanized habitats. This study investigates how human populations and urban development are associated with Northern Cardinal observation frequencies in U.S. counties. Given their adaptation to certain human-developed areas for food sourcing, this research examines the effect of urban development and population density on observation incidence. Focusing on 2022 data, the study analyses County-Level Rural-Urban Continuum Codes, Urban Influence Codes, Economic Typology, Natural Rate of Change, and Immigration Rates. The questions this research project aims to address are: What effect does urbanization have on Northern Cardinal sightings in the U.S.? It hypothesizes that more developed counties will record higher Northern Cardinal observations.

## **2. Datasets**

This project merged two datasets. The first dataset is from Cornell University’s Ornithology Lab’s eBird<sup>2</sup> project, providing Northern Cardinal observations by U.S. counties. Access was granted upon request, and data from January, April, July, and October 2022 were chosen to manage large file sizes and minimize seasonal biases, totalling 982,667 observations. Data cleaning involved aggregating data at the county level, retaining total observation counts, average observation duration, and effort distance. These variables, showing significant skewness, underwent log-transformation, including the key observation count variable for regression analysis. Appendix includes histograms of these variables before and after their log transformation (*Plots 1-3*). Nine states with minimal Northern Cardinal observations and counties with fewer than 10 observations in 2022 were dropped. Additionally, the top and bottom 1% of the observation counts were dropped, resulting in 2,288 aggregated counties in the cleaned dataset.

The second dataset, from the United States Department of Agriculture (USDA)<sup>3</sup>, contained county population and urbanization, including data on 3,143 counties. It required minimal cleaning, with no need for row eliminations due to missing values. The cleaning process primarily involved narrowing down to key columns: Rural Urban Continuum Code, Urban Influence, Economic Typology Code, Natural Rate of Change, and Net Immigration Rate.

The Rural-Urban Continuum Codes classify counties on a 1-9 scale, with Code 1 representing highly urbanized areas (over 1 million population) and Code 9 indicating rural areas (under 2,500 population, not metro-adjacent). Urban Influence Codes extend this categorization with a 1-12 range, where 1 is the most urban and 12 the most rural. More details of both codes will be found in *Table 1* and *Table 2*. Economic Typology Codes identify the primary economy of a county: 0 for non-specialized, 1 for Farm, 2 for Mining, 3 for Manufacturing, 4 for Federal/State Government, and 5 for Recreation Dependent. The net immigration rate measures population movement in and out of counties, while the natural rate of change measures the rate of births versus deaths. These rates, showing normal distribution, along with categorical development codes being converted to dummy variables, were used in regression models analysing Northern Cardinal observation counts. During the merging process, null and infinite

1. [https://www.allaboutbirds.org/guide/Northern\\_Cardinal/overview](https://www.allaboutbirds.org/guide/Northern_Cardinal/overview)

2. <https://ebird.org/data/download/ebd>

3. <https://www.ers.usda.gov/data-products/county-level-data-sets/county-level-data-sets-download-data/>

values were converted to 'NaN' (of which, there were below 10 missing values), and then filled with the mean value of their column. After merging the datasets, the study included 2,190 observations. The appendix provides detailed code descriptions distribution charts for the categorical county code variables and continuous variables (*Plots 4-8*). Note that distribution plots are located near the end of the appendix, as the regression models are considered of higher importance.

### 3. Regression Models

There were four sets of regression models run to test the association between urban development and populations on the log of Northern Cardinal observations. As the first three variables were categorical variables, they were run separately in groups of binary variables to avoid multicollinearity and overfitting by too many variables. A Variance Inflation Factor (VIF) test was run for each set of variables to identify potential multicollinearity, of which, none was found.

#### *3.1 OLS Regression Models of Urban Influence and Rural-Urban Continuum Codes on Log of Observation Count*

In this research, two distinct sets of Ordinary Least Squares (OLS) models were executed to examine the association of Rural-Urban Continuum Codes and Urban Influence Codes on the log of Northern Cardinal observations. The selection of the OLS methodology was based on the hypothesized linear association between the urban development status and the log of Northern Cardinal observations. Given that both Rural-Urban Continuum and Urban Influence Codes provide similar insights regarding a county's characteristics, they were analysed independently to enhance the robustness of the findings. In each set of models, the initial code (Code 1 for both sets) was the reference category and excluded from the regression analysis. Subsequently, a series of OLS regressions were conducted for each development code set, progressively incorporating additional development code variables as explanatory factors (Rural-Urban Continuum Codes 2 through 9 and Urban Influence Codes 2 through 12, respectively). To account for potential heteroscedasticity, all models were estimated using heteroscedasticity-consistent standard errors of type HC1. These models can be found in the appendix (*Regression Models 1 & 2*).

The analysis of both model sets yielded similar outcomes. In each case, the comprehensive model, including all codes except the benchmark (code 1), demonstrated superior model efficacy, as evidenced by the highest R-Squared value. These models were selected for analysis. The R-Squared for the Urban Influence Code OLS regression was determined to be 32.5%, while 34.1% for the Rural-Urban Continuum Code OLS regression. This indicates that approximately 33% of the variation in log observation counts is explainable by each model.

A consistent pattern emerged across both the Rural-Urban Continuum Codes and Urban Influence Codes. All codes exhibited statistically significant correlation coefficients at the 1% significance level. This uniform significance is likely due to the substantial data volume. Notably, in each model, the correlation coefficients increasingly deviated from zero in a negative trajectory. This trend suggests a positive relationship between the urbanization codes and the log of observation counts, which is observed in both the Rural-Urban Continuum Codes (Codes 2-9) and Urban Influence Codes (Codes 2-12). Therefore, the data provides evidence in support of a positive association between both Rural-Urban Continuum Codes and Urban Influence Codes, and the log of Northern Cardinal observations in US counties.

### **3.2 OLS Regression Models of County Economic Typology on Log Counts of Northern Cardinal Observations**

OLS models were executed to examine the association between County Typology Codes and the log of Northern Cardinal observations. The OLS model was selected due to the hypothesized linear association between the urban development status and the log of Northern Cardinal observation frequencies. Economic typology represents the urban development status of a county from a different. In the model, the code 0 (non-specialized counties) was designated as the reference category and excluded from the regression analysis. A series of OLS regressions were conducted for each typology code, progressively incorporating additional code variables as explanatory factors (Code 1: Economically Farming Dependent, etc.). To account for potential heteroscedasticity, all models were estimated using heteroscedasticity-consistent standard errors of type HC1. These models can be found in the appendix (*Regression Models 3*).

The analysis revealed that the regression model with all economic types was the best fit, with an R-Squared suggesting that 11.8% of the variance in log observation counts is explained by this model. Farming, mining, and manufacturing dependent resulted in statistically significant negative coefficients at the 1% threshold, while government dependent and recreation dependent were not statistically significant. This would suggest that farming, mining, or manufacturing dependent counties are associated with less Northern Cardinal observation counts.

### **3.3 OLS Regression Models for Natural Change Rate, Net Immigration Rate, Observation Duration, and Effort Distance on Log Observation Counts**

The final OLS models examined the associations between population changes and observation methods, with Northern Cardinal observation counts. These models incorporated natural change and net immigration rates per county to explore how population dynamics relate to observation counts. Additionally, the log-transformed variables for duration and distance of observations were included to examine potential associations with observation counts. The chosen model accounted for 20% of the variance in observation counts. Both natural change and net migration rates exhibited small, yet significant, positive associations with observation counts at the 1% significance level. For observation methods, the model included variables such as the log duration and log distance of observation. The log duration of observation displayed a positive association at the 1% significance threshold, indicating an association between longer observation times and higher counts. Conversely, the log distance travelled for observation had a negative association at the same significance threshold, indicating an association between greater travel distances for observations and lower counts. These regression models can be found in the appendix (*Regression Models 4*).

## **4. Generalization and external validity**

Predictive models were applied to the OLS models' variables to validate the associations with log Northern Cardinal observation counts. These models further validated the negative association between urban influence, rural-urban continuum codes, and observation counts, as lower observations counts were associated with rural counties (*Predictive Models 1 & 2*). Economic typology modelling indicated non-specialized counties and recreation dependent counties have the highest observation counts, with farming, mining, and manufacturing dependent counties showing lower counts (*Predictive Model 3*). Predictive models for natural change rate, net immigration rate, observation duration, and effort distance yielded results like the OLS models in their association with observation counts, with positive associations of natural change and net immigration rates seemingly stronger than in the regression models, a positive association with longer observation durations, and a negative association with greater

observation distances (*Predictive Models 4-7*). These findings suggest that observation counts of the Northern Cardinal are positively associated with urbanization levels, with recreational areas having strong associations. Observation duration positively correlates with counts, whereas wider observation effort distances may lead to fewer observations. This underscores the potential influence of observation methods and urban development on bird observation frequencies.

## 5. Causal Interpretation & Main Summary

Observation studies such as this cannot establish causality, but they can stimulate hypotheses for future research on the question. If in fact, urbanization does have a positive effect on the incidence of Northern Cardinal observation, it would lead to further research questions into what specific factors in the urbanized setting are causing this association. What other species of birds have similar associations with urbanization (e.g., ducks, geese, pigeons). Could the proportion of green space in a metropolitan setting influence this association? Could certainly urban settings demonstrate this association, while others do not? Does this association change based upon the industrial characteristics of the city?

The Urban Influence and Rural-Urban Continuum code models indicate a positive association between Northern Cardinal observations and urban development at the 1% significance threshold. This may reflect an attraction of these bird species to urbanized habitats, possibly due to food access and predator protection, despite initial concerns about observer bias in more populated areas. Economic typology models suggest a varying likelihood of observing backyard bird species across different urban habitats. Notably, recreation dependent counties, while not statistically significant, appear to host more of these species, contrasting with lower observation counts in farming, mining, and manufacturing dependent counties. This could be attributed to production sites destroying potential habitats for backyard bird species to reside, signifying the importance of environmental responsibility to preserve spaces for wildlife.

Regarding population dynamics, the statistically significant positive correlations of net immigration and natural change rates with bird observations suggest a potential link between higher human populations and more bird sightings. This might be due to easier food access for birds in human-dominated environments and increased likelihood of observation in more populated areas. Observation statistics further support these findings by showing longer observation durations positively correlate with higher bird counts, while greater travel distances during observations negatively impact counts. This aligns with the backyard habitat tendencies of Northern Cardinals, indicating that stationary observations in a single location, like a backyard, are more likely to yield higher observation counts.

## 6. Conclusion

The results of this study on Northern Cardinals suggests that the incidence of observing backyard birds in the U.S. is positively associated with urbanization at the county level. Three potential possibilities could explain this observation: 1. more observers in urbanized environments are documenting more observations, or 2. urbanization is responsible for higher populations of Northern Cardinals; or 3. both explanations are present.

Future research that includes population density datasets in multiple linear regression models analyzing the association between urbanization and Northern Cardinal observation could control population density, assuming multicollinearity is not present. Further research will also need to be done on the association between urban development and non-adapted bird species to analyze the impact on non-backyard bird species.

## 7. Appendix

### 7.1 Variable Description Tables

Table 1: Rural-Urban Continuum Codes for Counties Descriptions<sup>4</sup>

Code	Description
1	Metro areas with over 1M population
2	Metro areas with 250K-1M population.
3	Metro areas with under 250K population.
4	20K+ urban, metro-adjacent.
5	20K+ urban, non-metro-adjacent.
6	2.5K-20K urban, metro-adjacent.
7	2.5K-20K urban, non-metro-adjacent.
8	Rural or <2.5K urban, metro-adjacent.
9	Rural or <2.5K urban, non-metro-adjacent.

Table 2: Urban Influence Codes for Counties Descriptions<sup>5</sup>

Code	Description
1	Large metro area, 1M+ residents.
2	Small metro area, under 1M residents.
3	Micropolitan, adjacent to large metro.
4	Noncore, adjacent to large metro.
5	Micropolitan, adjacent to small metro.
6	Noncore, adjacent to small metro, town $\geq 2.5K$ residents.
7	Noncore, adjacent to small metro, no town $\geq 2.5K$ .
8	Micropolitan, not metro-adjacent.
9	Noncore, adjacent to micro area, town $\geq 2.5K$ .
10	Noncore, adjacent to micro area, no town $\geq 2.5K$ .
11	Noncore, not adjacent to metro/micro, town $\geq 2.5K$ .
12	Noncore, not adjacent to metro/micro, no town $\geq 2.5K$ .

4. <https://www.ers.usda.gov/data-products/rural-urban-continuum-codes/documentation/>

5. <https://www.ers.usda.gov/data-products/urban-influence-codes/documentation/>

## 7.2 Regression Models

### Regression Model 1: OLS Regression Models of Urban Influence Codes on Log of Observation Count

[illegible]

*Regression Model 2: OLS Regression Models of Rural-Urban Continuum Codes on Log of Observation Count*

	Dependent variable: <i>ln_OBSERVATION_COUNT</i>							
	Reg 12	Reg 13	Reg 14	Reg 15	Reg 16	Reg 17	Reg 18	Reg 19
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Rural-Urban Continuum Code 2	1.130*** (0.106)	1.257*** (0.108)	1.361*** (0.109)	1.363*** (0.110)	1.218*** (0.118)	0.847*** (0.127)	0.519*** (0.134)	-0.469*** (0.131)
Rural-Urban Continuum Code 3		0.898*** (0.102)	1.002*** (0.104)	1.004*** (0.104)	0.859*** (0.113)	0.488*** (0.122)	0.160 (0.130)	-0.829*** (0.127)
Rural-Urban Continuum Code 4			0.958*** (0.098)	0.961*** (0.099)	0.815*** (0.108)	0.444*** (0.117)	0.116 (0.125)	-0.872*** (0.122)
Rural-Urban Continuum Code 5				0.051 (0.174)	-0.095 (0.179)	-0.465** (0.185)	-0.793*** (0.191)	-1.782*** (0.189)
Rural-Urban Continuum Code 6					-0.419*** (0.088)	-0.789*** (0.100)	-1.117*** (0.109)	-2.106*** (0.105)
Rural-Urban Continuum Code 7						-1.295*** (0.112)	-1.623*** (0.121)	-2.611*** (0.118)
Rural-Urban Continuum Code 8							-1.528*** (0.143)	-2.517*** (0.140)
Rural-Urban Continuum Code 9								-3.025*** (0.125)
Constant	5.549*** (0.039)	5.422*** (0.042)	5.318*** (0.046)	5.316*** (0.047)	5.462*** (0.064)	5.832*** (0.079)	6.160*** (0.091)	7.149*** (0.086)
Observations	2190	2190	2190	2190	2190	2190	2190	2190
R <sup>2</sup>	0.049	0.077	0.098	0.098	0.106	0.152	0.190	0.341
Adjusted R <sup>2</sup>	0.049	0.076	0.097	0.097	0.104	0.150	0.187	0.339
Residual Std. Error	1.710 (df=2188)	1.685 (df=2187)	1.666 (df=2186)	1.667 (df=2185)	1.659 (df=2184)	1.617 (df=2183)	1.581 (df=2182)	1.426 (df=2181)
F Statistic	112.799*** (df=1; 2188)	92.499*** (df=2; 2187)	83.933*** (df=3; 2186)	63.098*** (df=4; 2185)	64.214*** (df=5; 2184)	85.025*** (df=6; 2183)	87.607*** (df=7; 2182)	145.267*** (df=8; 2181)
Note:								*p<0.1; **p<0.05; ***p<0.01

### Regression Model 3: OLS Regression Models of Economic Typology on Log of Observation Count

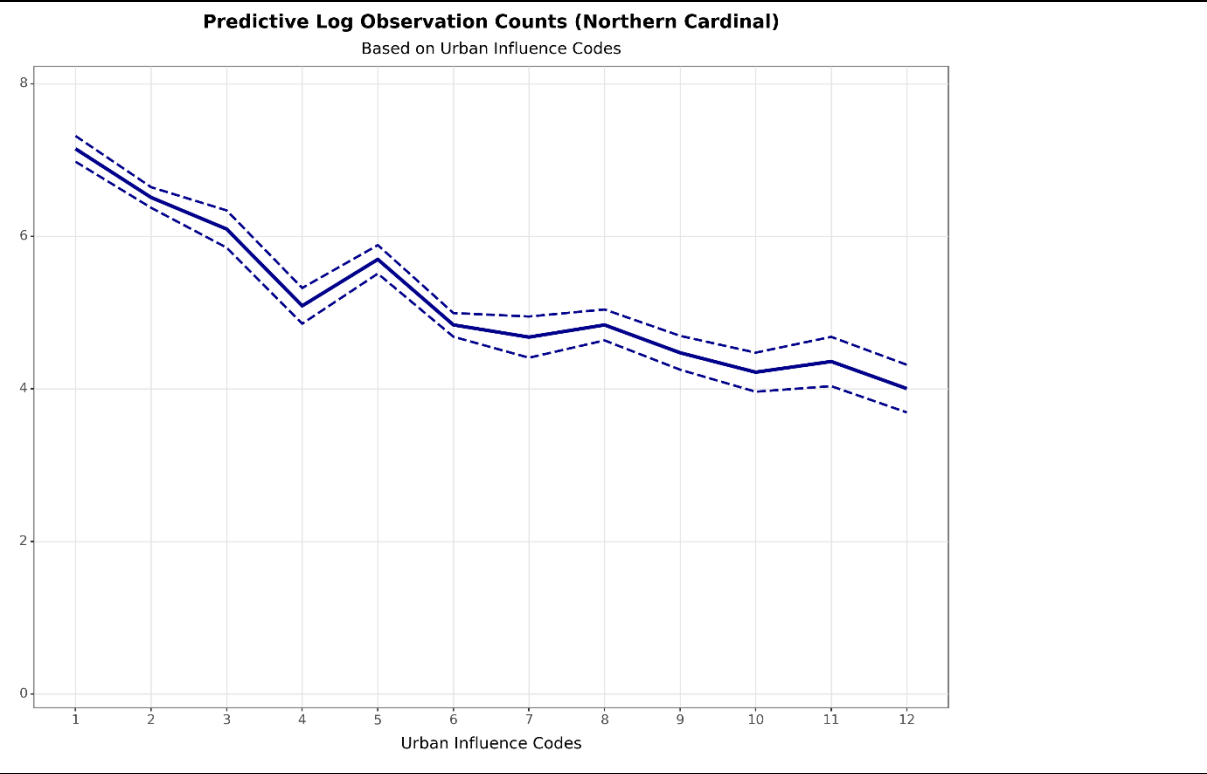
Dependent variable: ln_OBSERVATION_COUNT					
	Reg 20	Reg 21	Reg 22	Reg 23	Reg 24
	(1)	(2)	(3)	(4)	(5)
Economically Farming Dependent	-1.662*** (0.106)	-1.729*** (0.107)	-1.905*** (0.109)	-1.894*** (0.112)	-1.864*** (0.115)
Economically Mining Dependent		-1.074*** (0.141)	-1.250*** (0.143)	-1.239*** (0.145)	-1.209*** (0.147)
Economically Manufacturing Dependent			-0.754*** (0.083)	-0.742*** (0.086)	-0.713*** (0.090)
Economically Federal/State Government Dependent				0.058 (0.114)	0.088 (0.117)
Economically Recreation Dependent					0.166 (0.127)
Constant	5.843*** (0.038)	5.910*** (0.039)	6.086*** (0.046)	6.075*** (0.052)	6.045*** (0.058)
Observations	2190	2190	2190	2190	2190
R <sup>2</sup>	0.068	0.089	0.117	0.117	0.118
Adjusted R <sup>2</sup>	0.068	0.088	0.116	0.116	0.116
Residual Std. Error	1.693 (df=2188)	1.675 (df=2187)	1.649 (df=2186)	1.649 (df=2185)	1.649 (df=2184)
F Statistic	243.984*** (df=1; 2188)	148.919*** (df=2; 2187)	120.819*** (df=3; 2186)	90.798*** (df=4; 2185)	73.382*** (df=5; 2184)
Note:	*p<0.1; **p<0.05; ***p<0.01				

### Regression Model 4: OLS Regression Models for Natural Change Rate, Net Immigration Rate, Observation Duration, and Effort Distance on Log Observation Counts

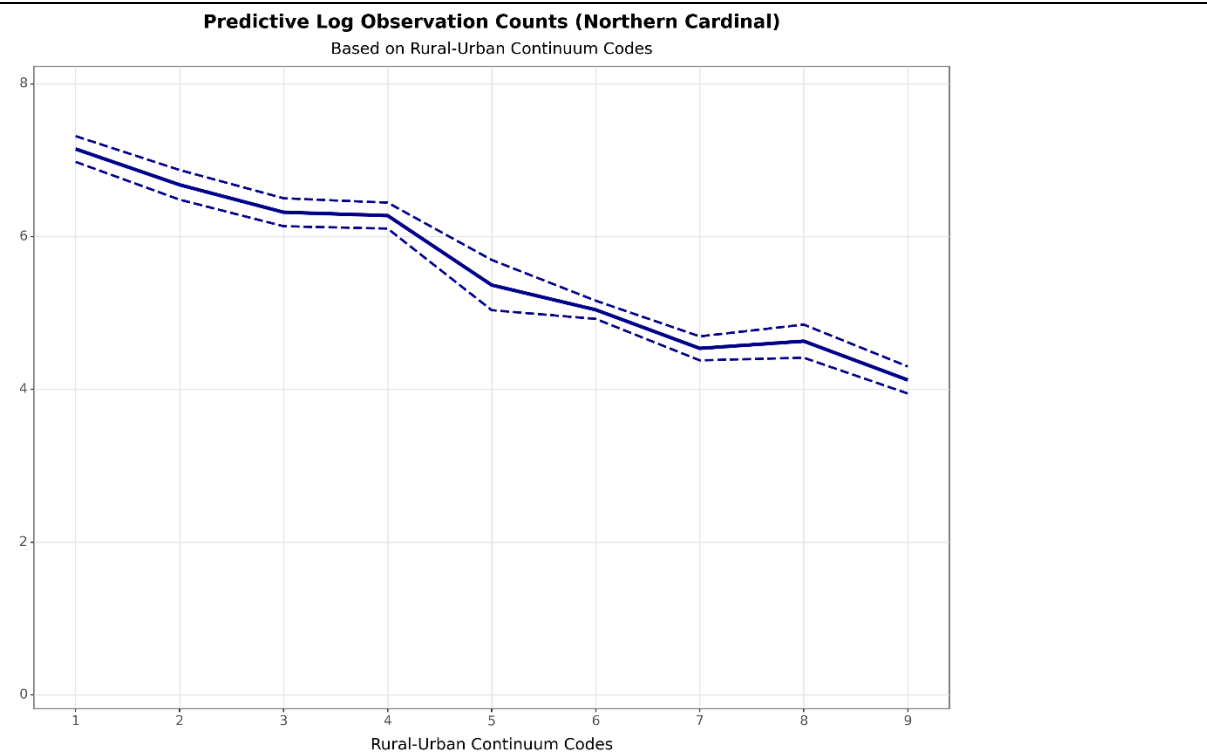
Dependent variable: ln_OBSERVATION_COUNT					
	Reg 25	Reg 26	Reg 27	Reg 28	Reg 29
	(1)	(2)	(3)	(4)	(5)
Natural Change Rate 2022	0.135*** (0.009)	0.141*** (0.009)		0.137*** (0.009)	0.133*** (0.009)
Net Migration Rate 2022		0.023*** (0.003)		0.021*** (0.003)	0.021*** (0.003)
Log Duration of Observation (minutes)			1.185*** (0.091)	0.866*** (0.079)	1.056*** (0.086)
Log Distance Traveled for Observation (km)			-0.569*** (0.072)		-0.462*** (0.069)
Constant	6.146*** (0.050)	6.023*** (0.053)	1.358*** (0.357)	2.384*** (0.332)	2.089*** (0.344)
Observations	2190	2190	2190	2190	2190
R <sup>2</sup>	0.106	0.135	0.085	0.183	0.200
Adjusted R <sup>2</sup>	0.106	0.134	0.085	0.182	0.199
Residual Std. Error	1.658 (df=2188)	1.632 (df=2187)	1.678 (df=2187)	1.586 (df=2186)	1.569 (df=2185)
F Statistic	220.342*** (df=1; 2188)	165.741*** (df=2; 2187)	90.868*** (df=2; 2187)	151.506*** (df=3; 2186)	125.026*** (df=4; 2185)
Note:	*p<0.1; **p<0.05; ***p<0.01				

7.3 Predictive Models

Predictive Model 1: Urban Influence Codes on Log Observation Counts



Predictive Model 2: Rural-Urban Continuum Codes on Log Observation Counts

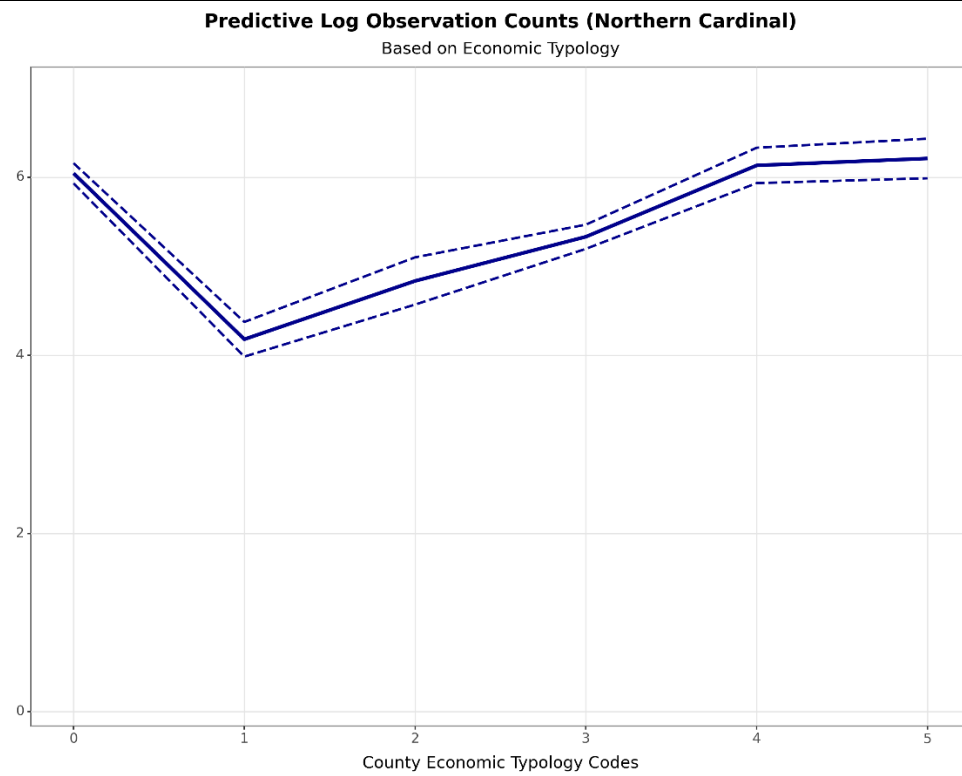




---

*Predictive Model 3: County Economic Typology Codes on Log Observation Counts*

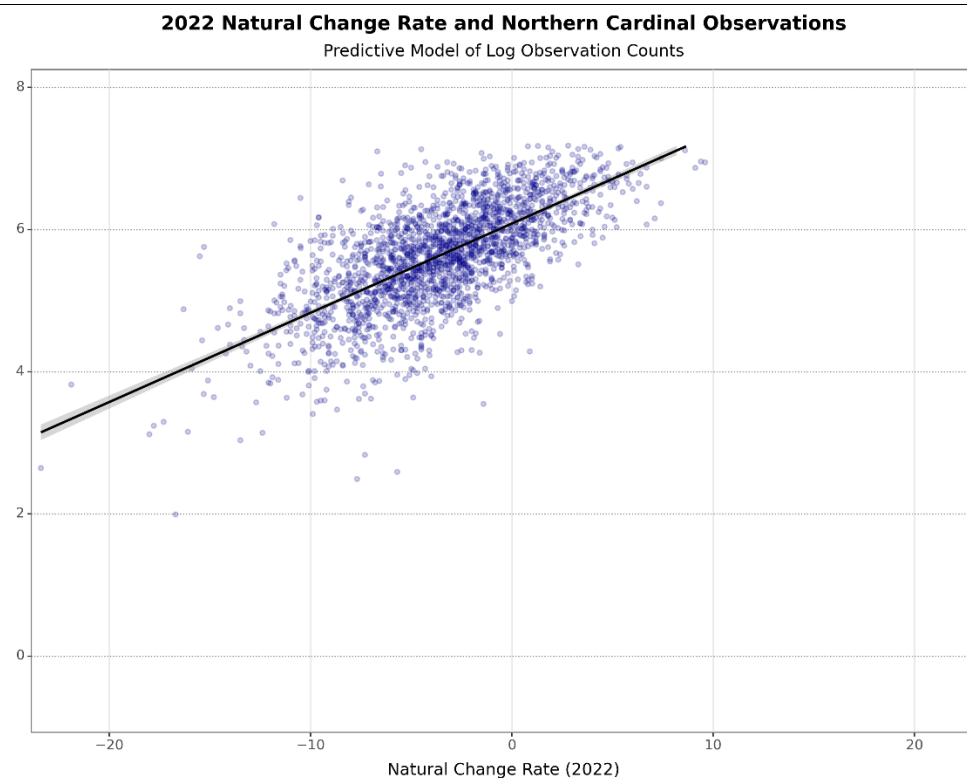
---



---

*Predictive Model 4: Natural Change Rate on Log Observation Counts*

---



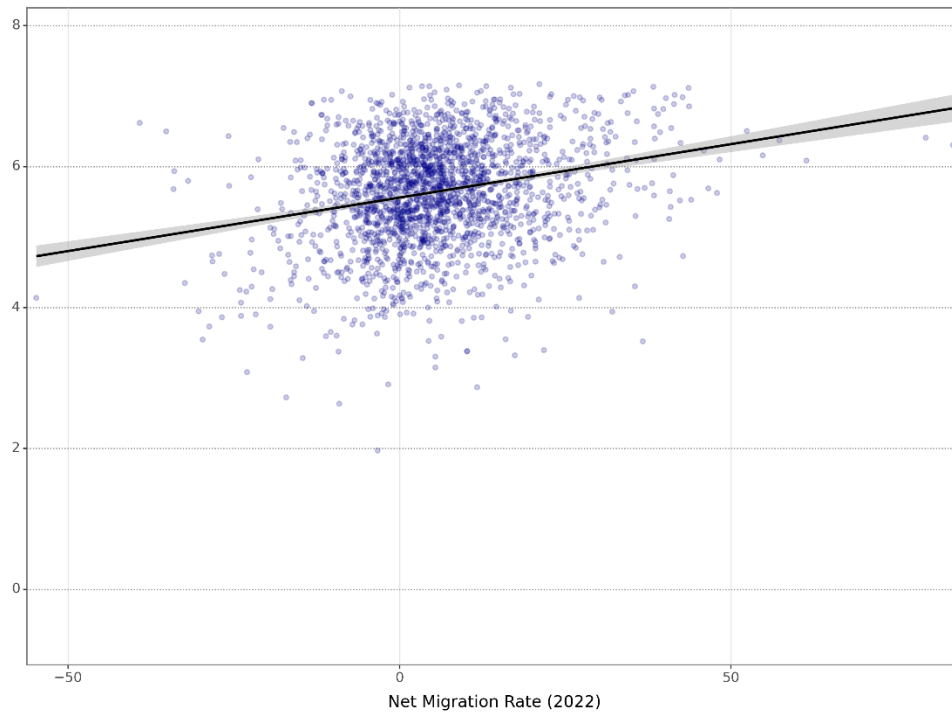
---

*Predictive Model 5: Net Immigration Rates on Log Observation Counts*

---

**2022 Net Migration Rate and Northern Cardinal Observations**

Predictive Model of Log Observation Counts



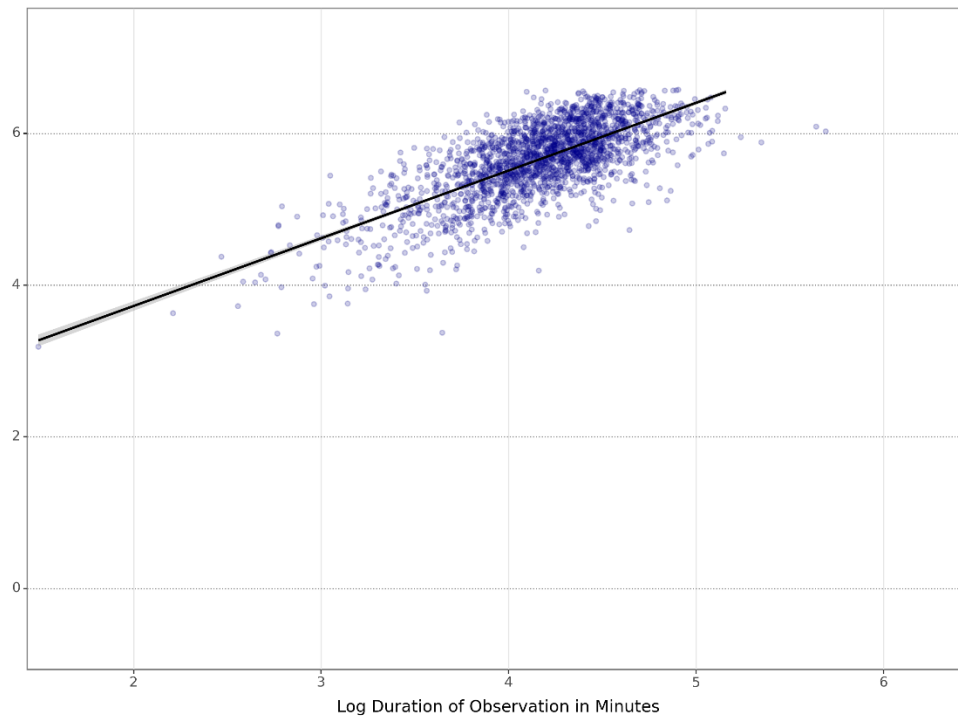
---

*Predictive Model 6: Duration of Observations on Log Observation Counts*

---

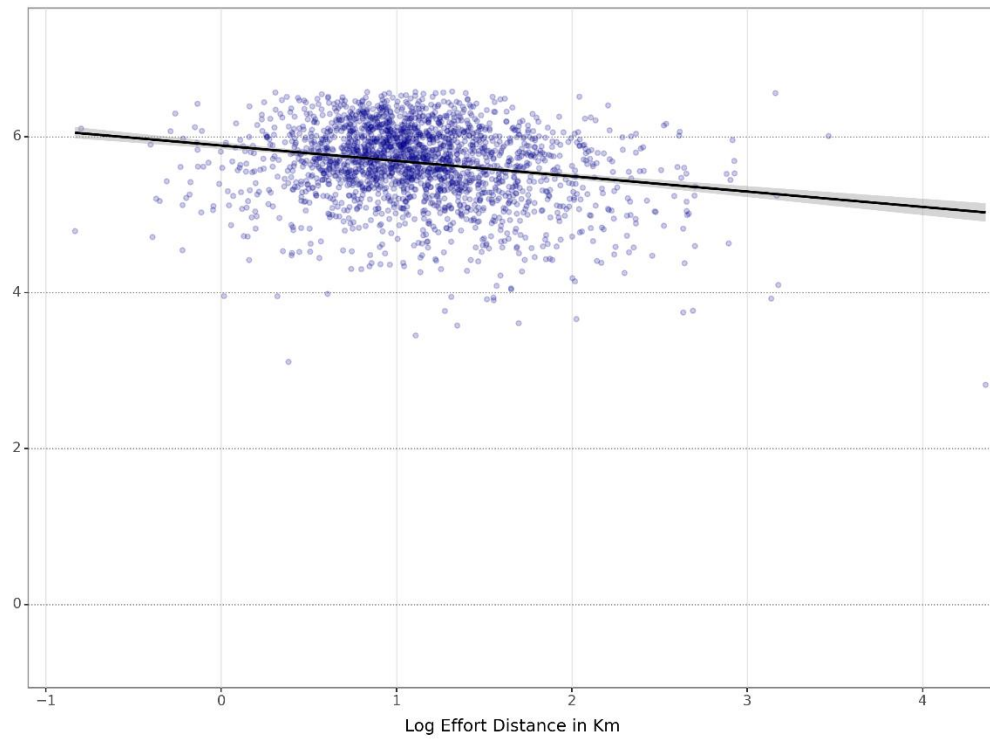
**Duration of Observations and Northern Cardinal Observation Counts**

Predictive Model of Log Observation Counts



*Predictive Model 7: Effort Distance on Log Observation Counts***Effort Distance (km) and Northern Cardinal Observation Counts**

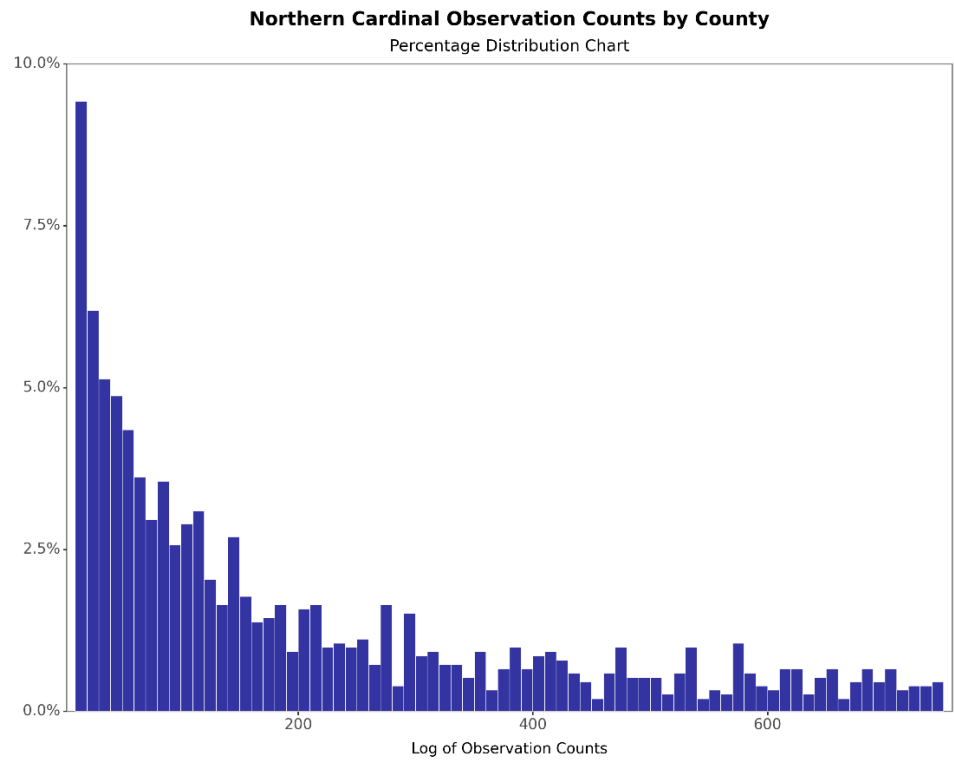
Predictive Model of Log Observation Counts



7.4 Distribution Charts

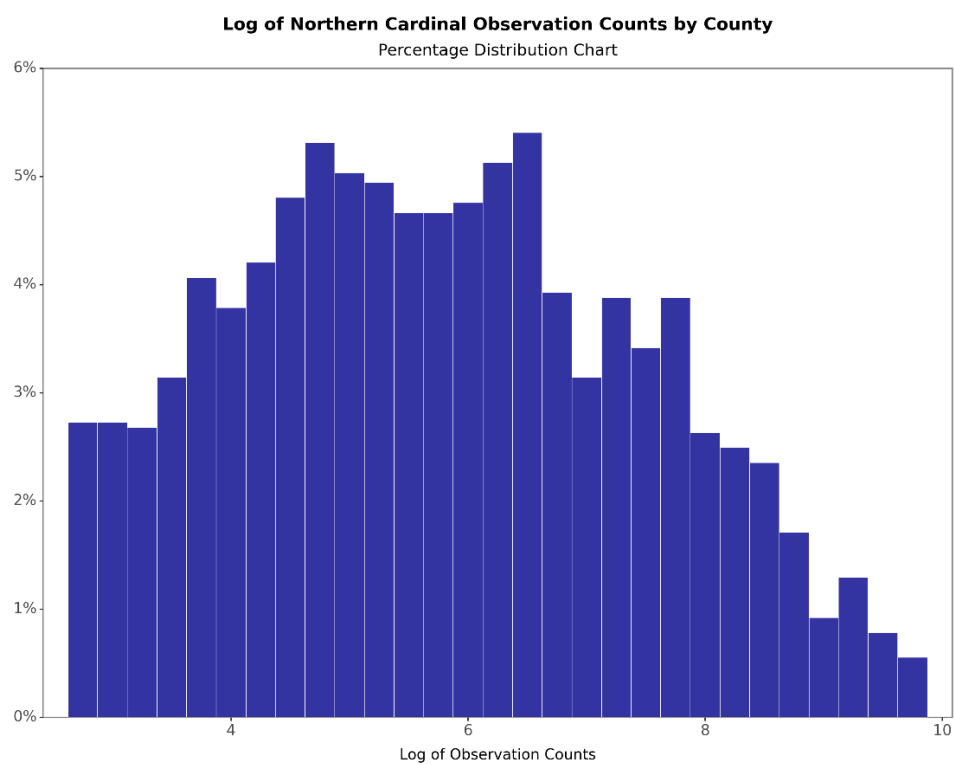
Plot 1(a & b): Distribution Plots of Observation Counts pre/post log-transformation.

a)



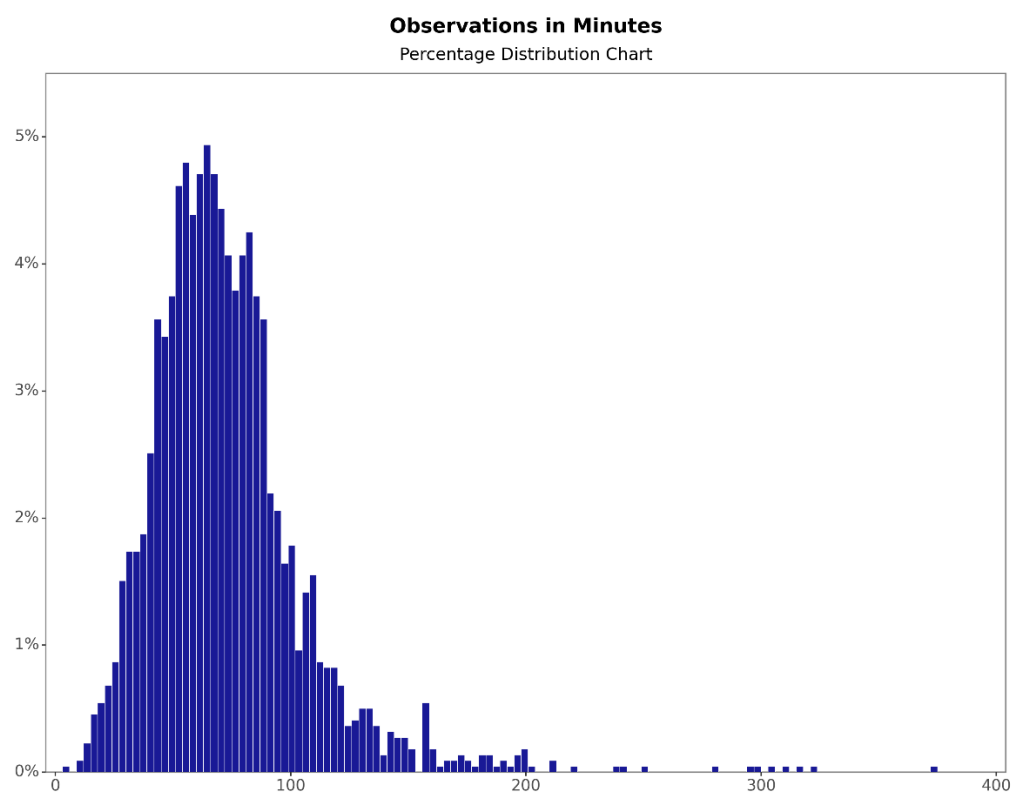
b)

---

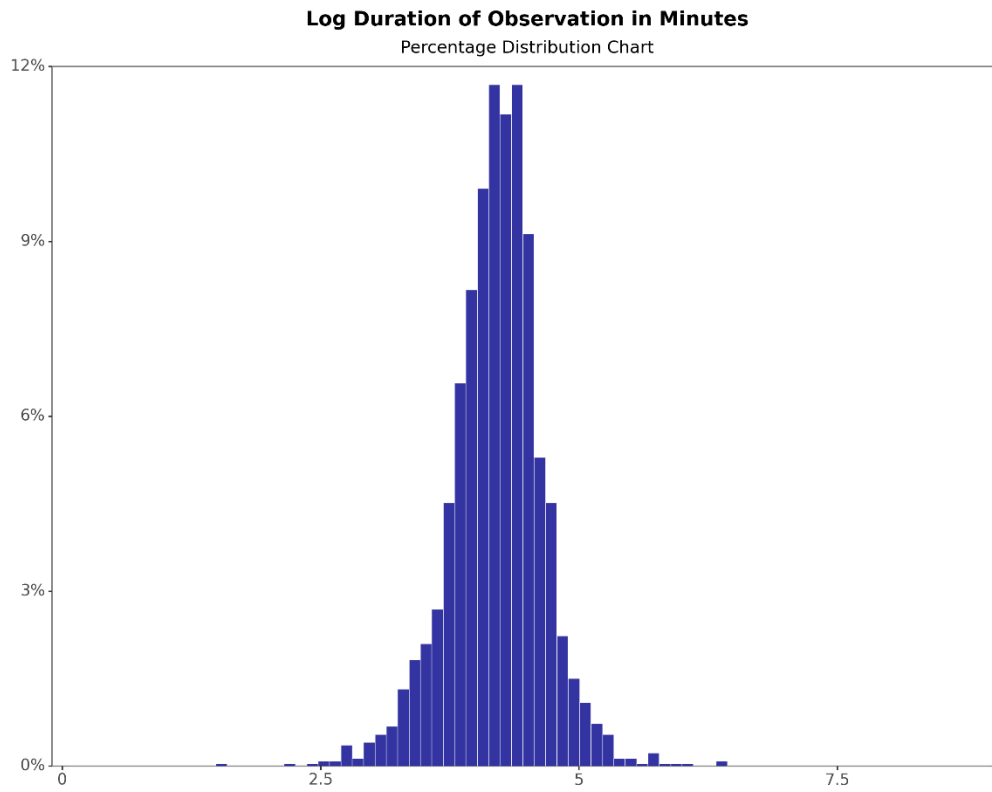


Plot 2 (a & b) Distribution Plots of Observation Duration pre/post log-transformation

a)



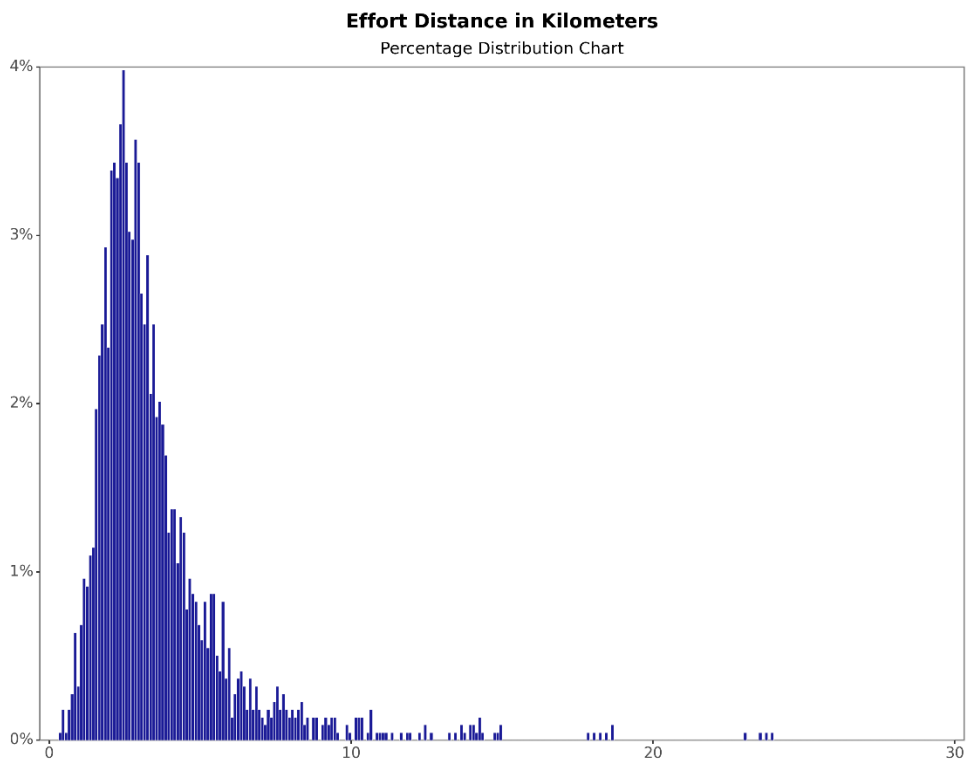
b)



Plot 3(a & b): Distribution Plots of Observation Distance pre/post log-transformation

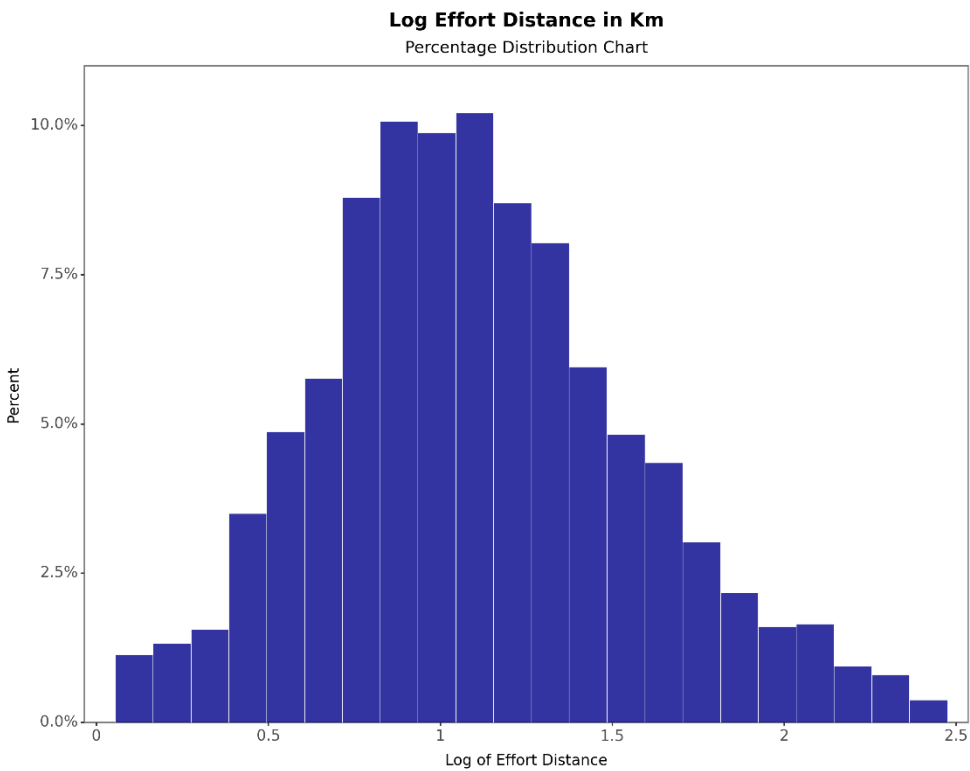
a)

---



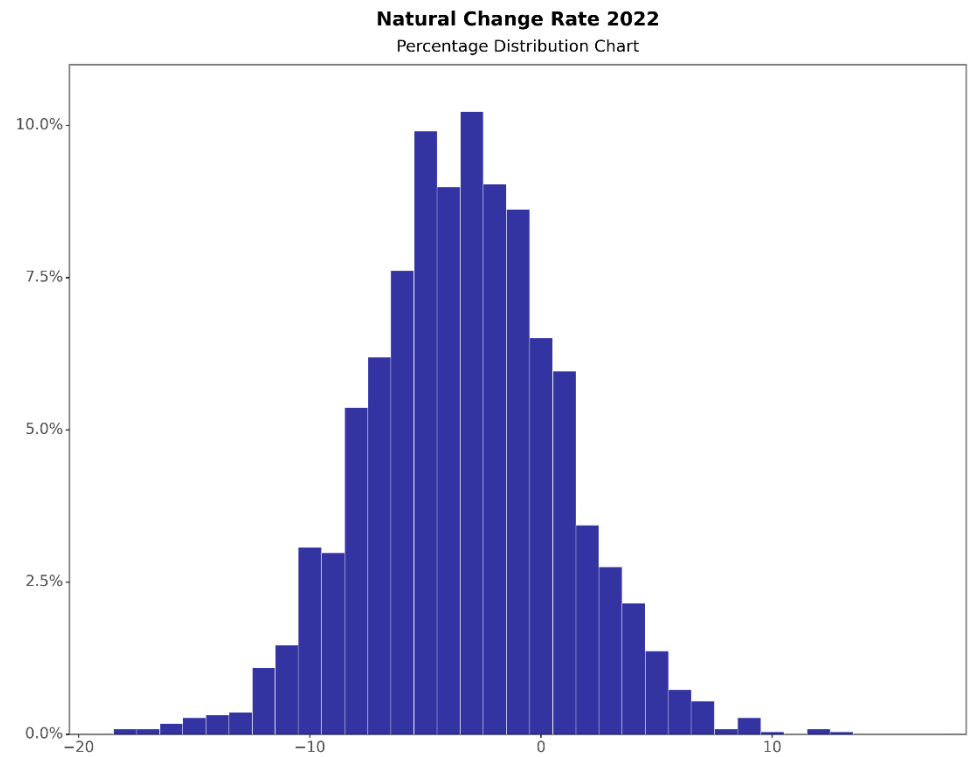
b)

---



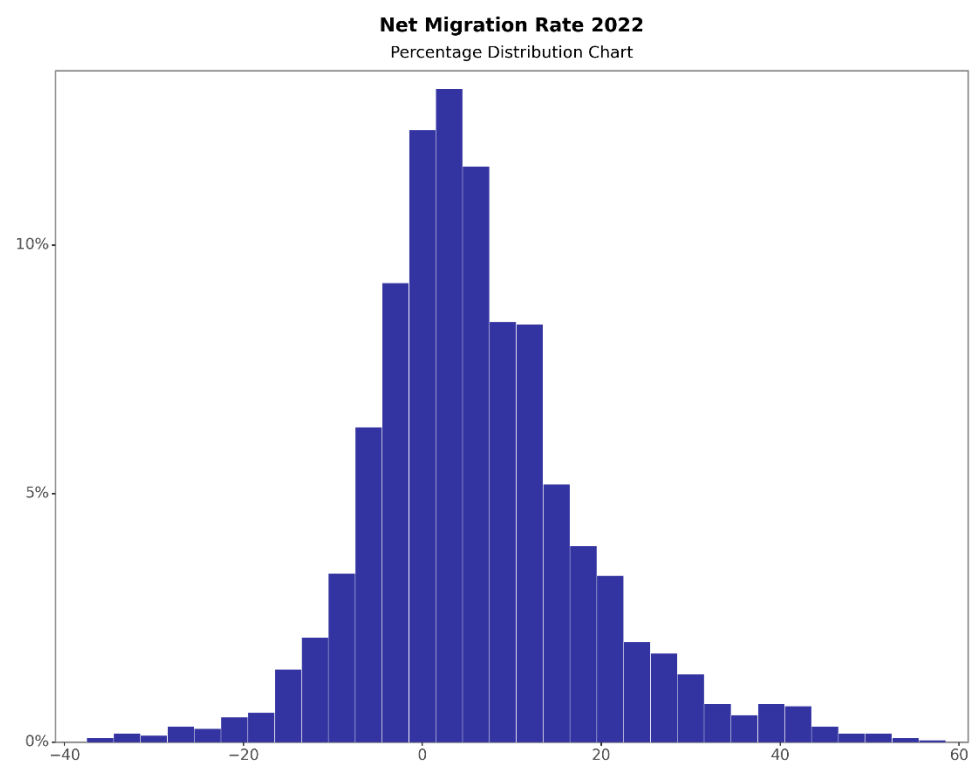
Plot 4: Distribution of Natural Change Rate in 2022

---

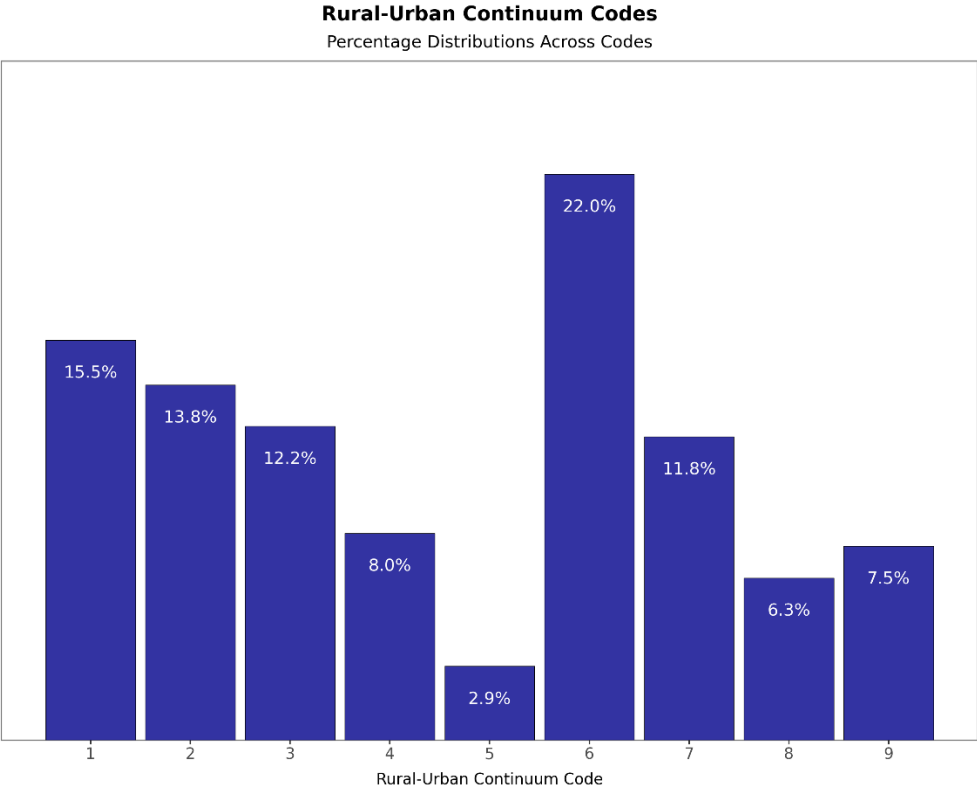


Plot 5: Distribution of Net Migration Rate 2022

---



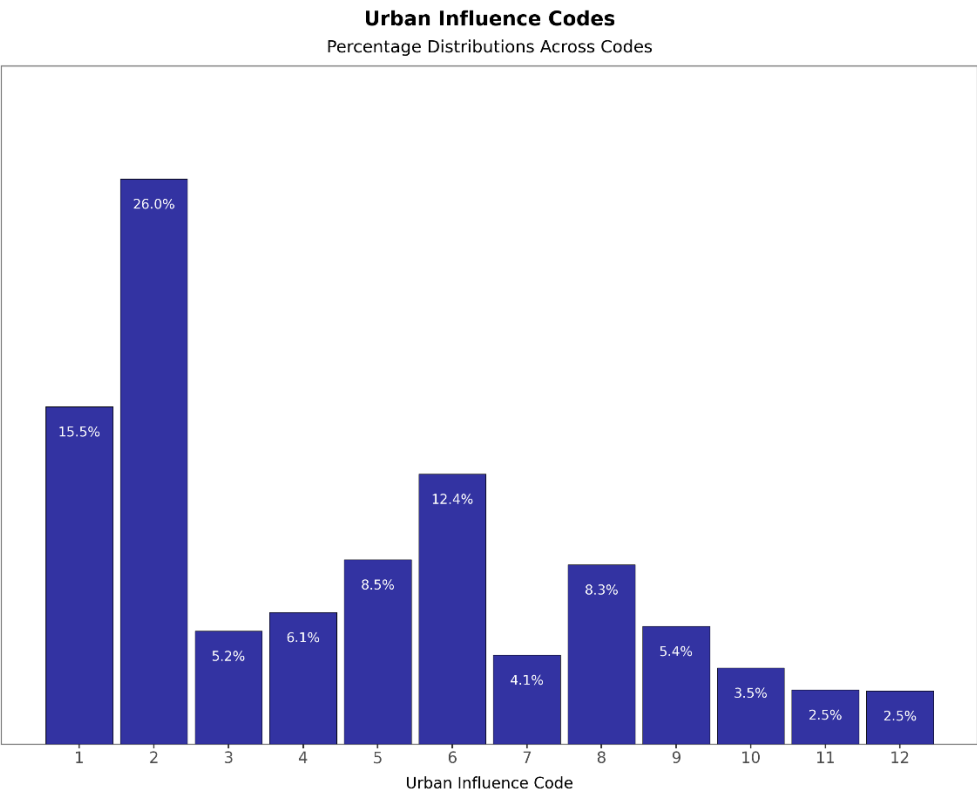
Plot 6: Distribution of Rural-Urban Continuum Codes





Plot 7: Distribution of Urban Influence Codes

---



Plot 8: Distribution of County Economic Typology in 2015

---

