

Murder: Analysis and Description

Problem Definition	Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)
--------------------	---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------



Murder: Analysis and Description

Problem Definition	Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)
--------------------	---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------

Primary response variable: Homicide Rate

This variable demonstrates homicide rates per 100,000 people per state. State is the identifying variable.

Classification variable: High/Low Homicide

This variable is classified as a “High” or a “Low” based on whether the homicide rate is above or below the national average. In some analyses 1 identifies a “High” and 0 identifies a “Low”.

Murder: Analysis and Description

Problem Definition	Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)
--------------------	---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------

Abortion Ban
This variable shows the state’s abortion ban by latest possible week of abortion. Assuming an average length of pregnancy of 40 weeks (9 months), the ban number (22, 26, or 40) depicts the last week in which a woman may get an abortion. If a state’s value is 40 weeks, there is no ban.

Prison Rate
This is the number of incarcerated people per 100,000 people per state. Prisons differ from jails in that prisons are longer-term incarceration facilities to which prisoners are moved following their sentencing. In general, while people wait for trials before they’re sentenced, they’re put in jail. After they’re sentenced (*“You are hereby sentenced to 12 months…”*), they go to prison.

Population Density
This value is calculated based on the total population of a state divided by the total land area. Higher values in this variable indicate that the state has a high number of residents who therefore live relatively close together.

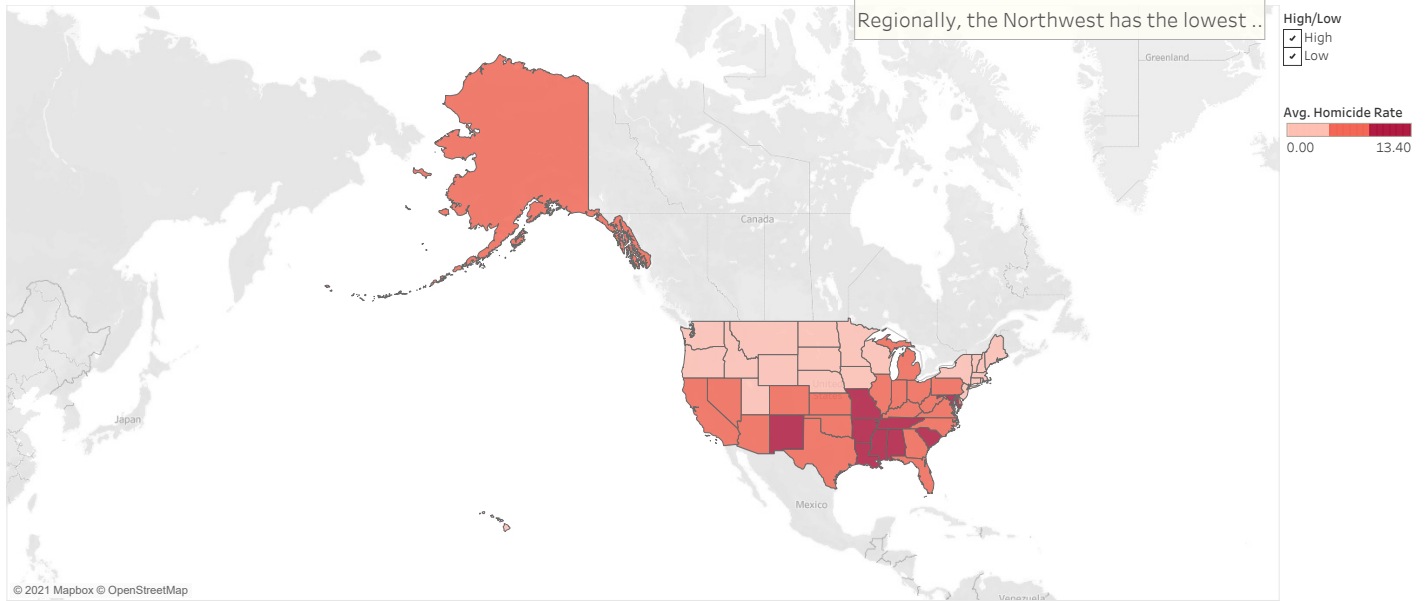
Governor Party
The governor party (Democratic, Independent, or Republican) is determined based on the predominant state governor from 2018. If a gubernatorial election was held during 2018 and a former governor was unseated, the governor that spent the majority of 2018—that is, at least 183 days—in the governor’s seat was considered the “governor of 2018.” There was only one Independent governor that year, and he was from Alaska.

Minority Population
Minority population is comprised of the groups “Black or African American”, “American Indian and Alaska Native”, and “Asian”. There are only a few states with a minority-majority population (for example, Hawaii has more Asian-identifying than White-identifying people), therefore, this dataset considers the largest nonwhite population in the state.

Region

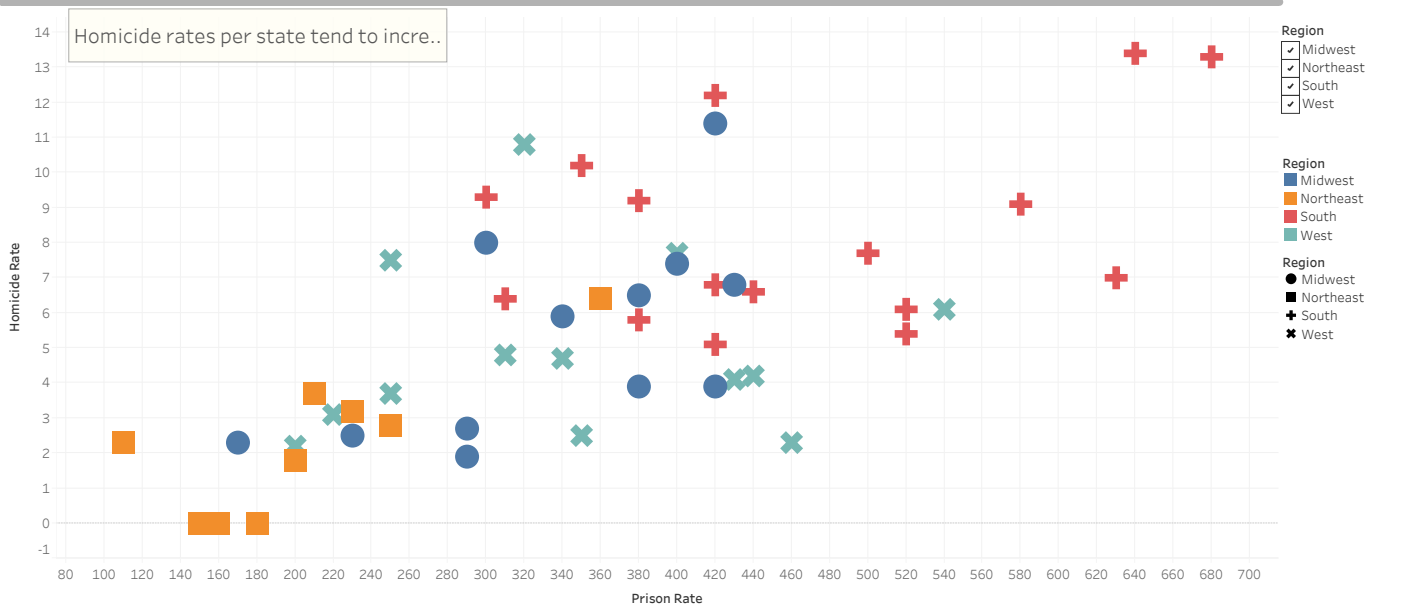
Murder: Analysis and Description

Problem Definition	Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)
--------------------	---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------



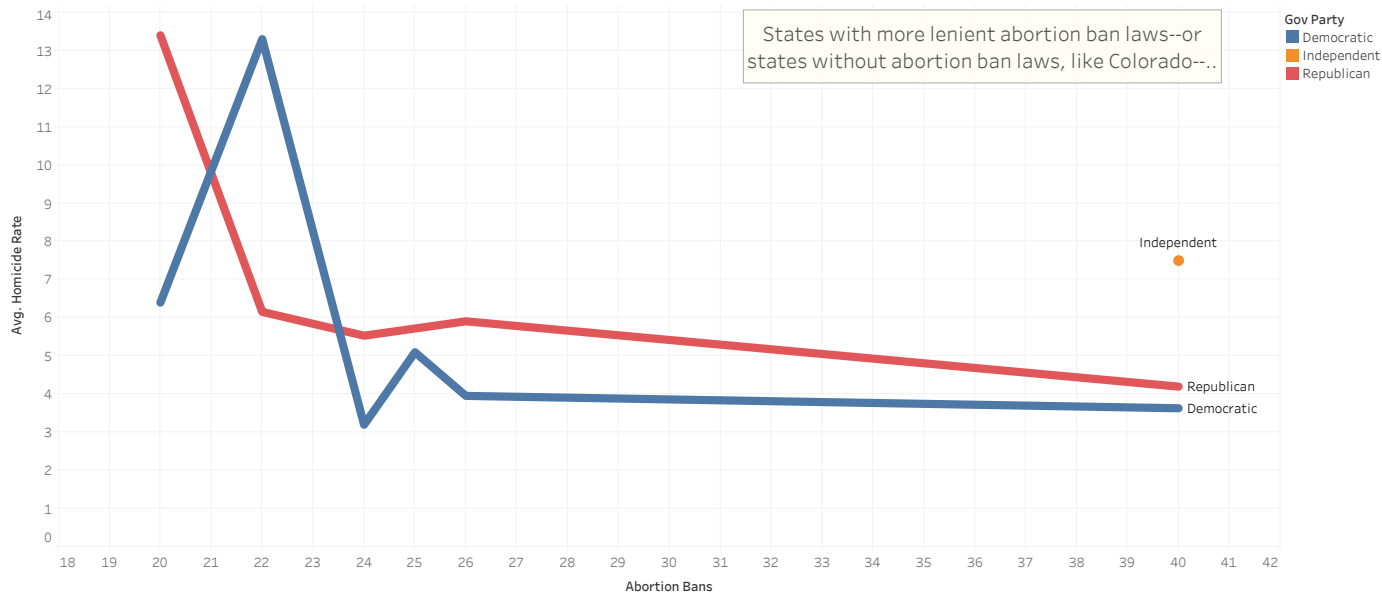
Murder: Analysis and Description

Problem Definition	Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)
--------------------	---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------



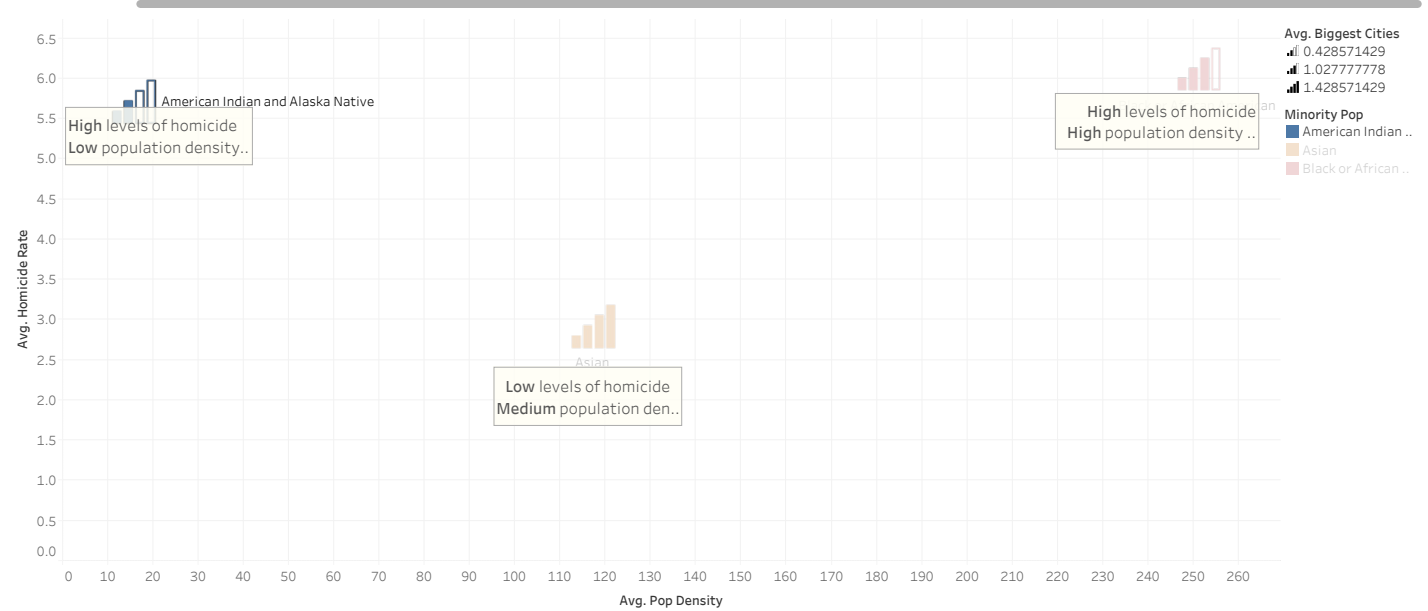
Murder: Analysis and Description

Problem Definition	Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)
--------------------	---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------



Murder: Analysis and Description

Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)	Collection
---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------	------------



Murder: Analysis and Description

Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)	Collection
---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------	------------

Continuous

Continuous

Model Comparison

Predictors

Measures of Fit for Homicide_Rate

Predictor	Creator	.2	.4	.6	.8	RSquare	RASE	AAE	Freq
M4: Predicted Values	Fit Least Squares					0.5888	2.1216	1.7414	50
M5: Predicted Values	Fit Least Squares					0.5337	2.2593	1.9118	50
M6: Predicted Values	Fit Least Squares					0.5519	2.2148	1.9070	50

Here we can see the comparison of each of the Continuous models:

M4: Least Squares

Summary of Fit

RSquare	0.588753
RSquare Adj	0.489708
Root Mean Square Error	2.433685
Mean of Response	5.614
Observations (or Sum Wgts)	50

the model with the best RSquare value is the Fit of Least Squares (Model 4) with a value of 58.9% which means the model can accurately predict 58.9% of the data based on the inputted values.

M5: Backward Stepwise

Summary of Fit

RSquare	0.533663
RSquare Adj	0.492211
Root Mean Square Error	2.381489
Mean of Response	5.614
Observations (or Sum Wgts)	50

However, based on the adjusted RSquare, the best model is the Forward stepwise sitting at 50.1%.

M6: Forward Stepwise

Summary of Fit

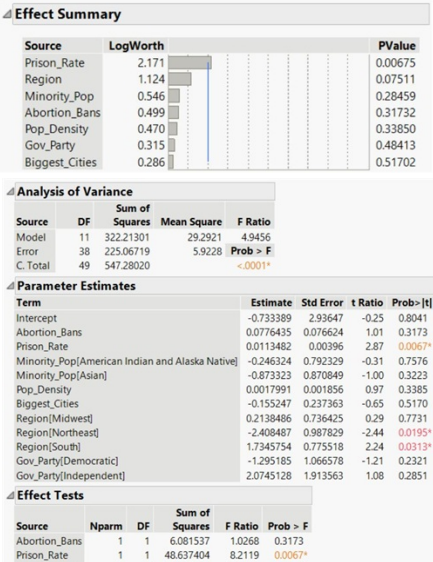
RSquare	0.551858
RSquare Adj	0.500933
Root Mean Square Error	2.360947
Mean of Response	5.614
Observations (or Sum Wgts)	50

Further analysis on this aspect is and model is in Continuous Analysis (pt.2)

Murder: Analysis and Description

Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)	Collection
---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------	------------

Continuous pt2



After selecting the best model of continuous data (Fit of Least Squares Regression) we can determine what variables are the best at predicting homicide rates.

Here we can see that only one variable passes the p-test of being less than .05 and that is Prison_Rate. This means that Prison_Rate is significant and is a better predictor of homicide rate than other variables. The next two values (variables) are Region[NorthEast] and Region[South] that are also decent but not the best predictors of homicide rate. Together the equation for this regression model is determined through the Parameter Estimates.

Murder: Analysis and Description

Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)	Collection
---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------	------------

Categorical

Model 1: Classification Tree

		M1: Most Likely High/Low	
High/Low	Validation	0	1
0	Training	13	2
	Validation	8	2
1	Training	1	14
	Validation	2	8

Training Error: $3/(13+2+1+14) = .1$
Validation Error: $4/20 = .2$
Training Sensitivity: $14/15 = .933$
Validation Sensitivity: $8/10 = .8$

Model 2: kNN

		M2: Most Likely High/Low	
High/Low	Validation	0	1
0	Training	13	2
	Validation	8	2
1	Training	1	14
	Validation	2	8

Training Error: $3/(13+2+1+14) = .1$
Validation Error: $4/20 = .2$
Training Sensitivity: $14/15 = .933$
Validation Sensitivity: $8/10 = .80$

Model 3: Naive Bayes

		Training		Validation	
		Actual	Predicted Count	Actual	Predicted Count
High/Low	0	0	1	0	1
	1	15	0	9	1
	1	0	15	2	8

Training Error: $0/15 = 0$
Validation Error: $3/20 = .15$
Training Sensitivity: $15/15 = 1$
Validation Sensitivity: $8/10 = .8$

Here we can see the analysis of the Categorical models:

Throughout these models, when conducting the error and sensitivity analysis, one model stands out: Naive Bayes.

With a training error of 0%, a validation error of 15% & training sensitivity of 100% and validation sensitivity at 80%, this model is nearly perfect. The lower the error rate and the higher the sensitivity rate, the better the model. Here, Naive Bayes can correctly predict whether homicide rate will be high or low with 80% certainty.

Murder: Analysis and Description

Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)	Collection
---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------	------------

Training				Validation			
Misclassification		Rate	Misclassifications	Misclassification		Rate	Misclassifications
Count				Count			
30		0	0	20		0.15000	3

Confusion Matrix			
Training			
Actual High/Low	Predicted Count		
	0	1	
0	15	0	
1	0	15	

Validation			
Actual High/Low	Predicted Count		
	0	1	
0	9	1	
1	2	8	

Fit Details			
Measure	Training	Validation	Definition
Entropy RSquare	0.9657	-0.090	$1 - \frac{\text{Loglike}(\text{model})}{\text{Loglike}(0)}$
Generalized RSquare	0.9838	-0.177	$\frac{(1 - (L(0)/L(\text{model}))^{(2/n)})}{(1 - L(0)^{(2/n)})}$
Mean -Log p	0.0238	0.7555	$\sum -\text{Log}(p_{ij})/n$
RASE	0.0651	0.3763	$\sqrt{\sum (y_{ij} - p_{ij})^2/n}$
Mean Abs Dev	0.0212	0.1847	$\sum y_{ij} - p_{ij} /n$
Misclassification Rate	0.0000	0.1500	$\sum (p_{ij} \neq p_{\text{Max}})/n$
N	30	20	n

As stated previously, there were 0 misclassifications within the training data set and only 3 misclassifications within the validation set. This shows the validity of Naive bayes theorem for this data set.

The Sensitivity Rate shows that the model can accurately predict 80% (validation) of the homicide rates whether they are / were high or low. This is very valuable to take into account for predicting future homicide rates based off of previous values.

Murder: Analysis and Description

Dependent Variables	Independent Variables	Visualization: Map	Visualization: Scatterplot	Visualization: Line Graph	Visualization: Point Graph	Continuous Analysis (pt.1)	Continuous Analysis (pt.2)	Categorical Analysis (pt.1)	Categorical Analysis (pt.2)	Collection
---------------------	-----------------------	--------------------	----------------------------	---------------------------	----------------------------	----------------------------	----------------------------	-----------------------------	-----------------------------	------------

Citations:

Ballotpedia (2018). Partisan Composition of Governors by State. Ballotpedia. Retrieved from https://ballotpedia.org/Partisan_composition_of_governors

Center for Disease Control (2018). Homicide Mortality by State. National Center for Health Statistics . Retrieved from https://www.cdc.gov/nchs/pressroom/sosmap/homicide_mortality/homicide.htm

Guttmacher Institute (2020). An Overview of Abortion Laws. State Laws and Policies . Retrieved from <https://www.guttmacher.org/state-policy/explore/overview-abortion-laws>

National Geographic (2021). United States Regions. National Geographic . Retrieved from <https://www.nationalgeographic.org/maps/united-states-regions/>

States101 (2021). U.S. States: Populations, Land Area, and Density. States101 . Retrieved from <https://www.states101.com/populations>

White, M. (2019). The Top 10 Largest U.S. Cities by Population. Moving.com . Retrieved from <https://www.moving.com/top-10-largest-cities-by-population/>