# Project 3: Thematic Mapping & Clustering
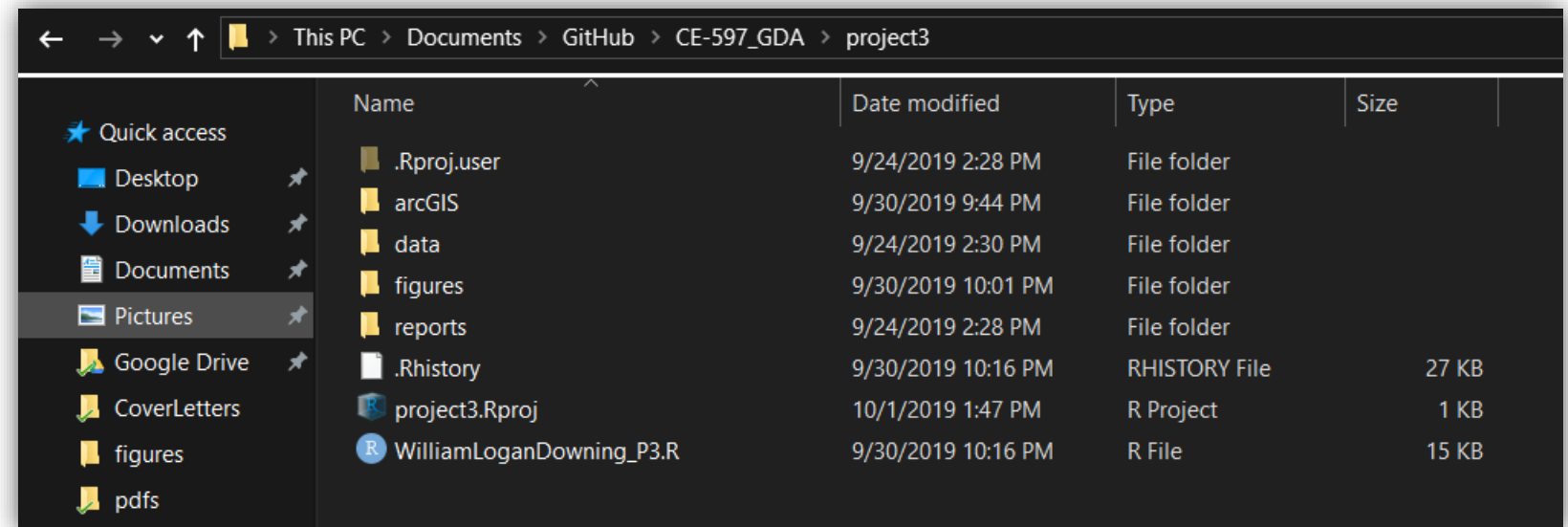
## W. Logan Downing

# Outline

1. Data Preparation
2. Thematic Mapping of Population Density
3. Develop a Kmeans Clustering Algorithm
4. Kmeans Comparison with ArcGIS

# Note

- You will see me using the png() command. This is to save my figures and persist them in the project structure that I've been using to keep myself organized.

An R Project is used to make R Studio use the source file directory as its working directory.

# Data Preparation

1. In this section, I opted to work with the tracts data **but** I have also chosen to read in and subset West Lafayette from the tracts for practice.

```
16  # Load the block data and clip it first.
17  blocks <- sf::st_read('./data/raw/tabblock2010_18_pophu', layer='tabblock2010_18_pophu')
18  #blocks_proj <- sf::st_transform(blocks, crs=32616)
19  # subset blocks to west lafayette.
20  wl <- osmdata::getbb('west lafayette, usa')
21  # need to create a polygon from the bbox.
22  coords <- matrix(c(wl[1,1], wl[2,1], wl[1,1], wl[2,2], wl[1,2], wl[2,2],
23                     wl[1,2], wl[2,1], wl[1,1], wl[2,1]), ncol=2,
24                   byrow=T)
25  geoobj <- st_sfc(st_polygon(list(coords)))
26  wl_poly <- st_sf(geoobj, crs=st_crs(blocks))
27
28  wl_blocks <- st_intersection(blocks, wl_poly)
29  wl_blocks <- sf::st_transform(wl_blocks, crs=32616)
30  sf::st_crs(wl_blocks)
31  wl_blocks$blockArea <- sf::st_area(wl_blocks) %>% set_units(km^2) # set the blockArea
```

PURDUE
UNIVERSITY

# Data Preparation

```
34 tract <- sf::st_read('./data/raw/Tract_2010Census_DP1_IN/utm', layer='tract_in_selected_utm')
35 sf::st_crs(tract)
36 # epsg code found at https://spatialreference.org/ref/epsg/wgs-84-utm-zone-16n/
37 tract_proj <- sf::st_transform(tract, crs=32616)
38
39 tract_proj$blockArea <- sf::st_area(tract_proj) %>% set_units(km^2)
40 tract_proj <- tract_proj[,c("GEOID10", "NAMELSAD10", "ALAND10", "AWATER10",
41                              "INTPTLAT10", "INTPTLON10", "DP0010001", "blockArea")]
42 class(tract_proj)
```

```
> tract <- sf::st_read('./data/raw/Tract_2010Census_DP1_IN/utm', layer='tract_in_selected_utm')
Reading layer `tract_in_selected_utm' from data source `C:\Users\Downi\Documents\GitHub\CE-597_GDA\pro
ject3\data\raw\Tract_2010Census_DP1_IN\utm' using driver `ESRI Shapefile'
Simple feature collection with 1511 features and 194 fields
geometry type:  POLYGON
dimension:      XY
bbox:           xmin: 403477.8 ymin: 4180915 xmax: 692181.8 ymax: 4625475
epsg (SRID):    26916
proj4string:    +proj=utm +zone=16 +datum=NAD83 +units=m +no_defs
> sf::st_crs(tract)
Coordinate Reference System:
  EPSG: 26916
  proj4string: "+proj=utm +zone=16 +datum=NAD83 +units=m +no_defs"
> # epsg code found at https://spatialreference.org/ref/epsg/wgs-84-utm-zone-16n/
> tract_proj <- sf::st_transform(tract, crs=32616)
>
> tract_proj$blockArea <- sf::st_area(tract_proj) %>% set_units(km^2)
> tract_proj <- tract_proj[,c("GEOID10", "NAMELSAD10", "ALAND10", "AWATER10",
+                              "INTPTLAT10", "INTPTLON10", "DP0010001", "blockArea")]
> class(tract_proj)
[1] "sf"         "data.frame"
```

PURDUE
UNIVERSITY

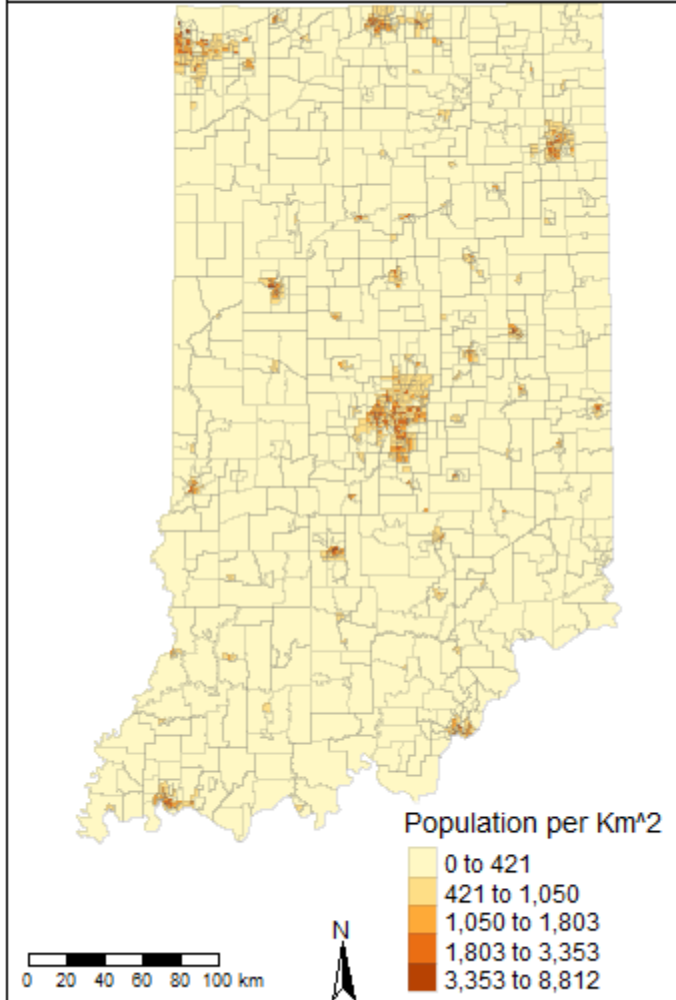# Thematic Mapping of Population Density

1. For this section, we needed to create thematic maps from the tract population (DP0010001).

```
67  k <- tm_shape(tract_proj) + tm_polygons('DP0010001', style='kmeans', title='Block Area (Km^2)') +
68     tm_compass(position=c('right', 'bottom'), text.size=1) +
69     tm_scale_bar(position=c('left', 'bottom'), text.size=1.5) +
70     tm_layout(frame=T, panel.show=T, main.title.position='center', main.title.size=1,
71              legend.position=c('right', 'bottom'), legend.text.size=1,
72              legend.title.size=1.4, legend.bg.color='white', legend.bg.alpha=.6,
73              inner.margins = c(.17, 0, 0, 0),
74              panel.labels='Population Clustered by Kmeans')
75
76  j <- tm_shape(tract_proj) + tm_polygons('DP0010001', style='jenks', title='Block Area (Km^2)') +
77     tm_compass(position=c('right', 'bottom'), text.size=1) +
78     tm_scale_bar(position=c('left', 'bottom'), text.size=1.5) +
79     tm_layout(frame=T, panel.show=T, main.title.position='center', main.title.size=1,
80              legend.position=c('right', 'bottom'), legend.text.size=1,
81              legend.title.size=1.4, legend.bg.color='white', legend.bg.alpha=.6,
82              inner.margins = c(.17, 0, 0, 0),
83              panel.labels='Population Clustered by Jenks')
```
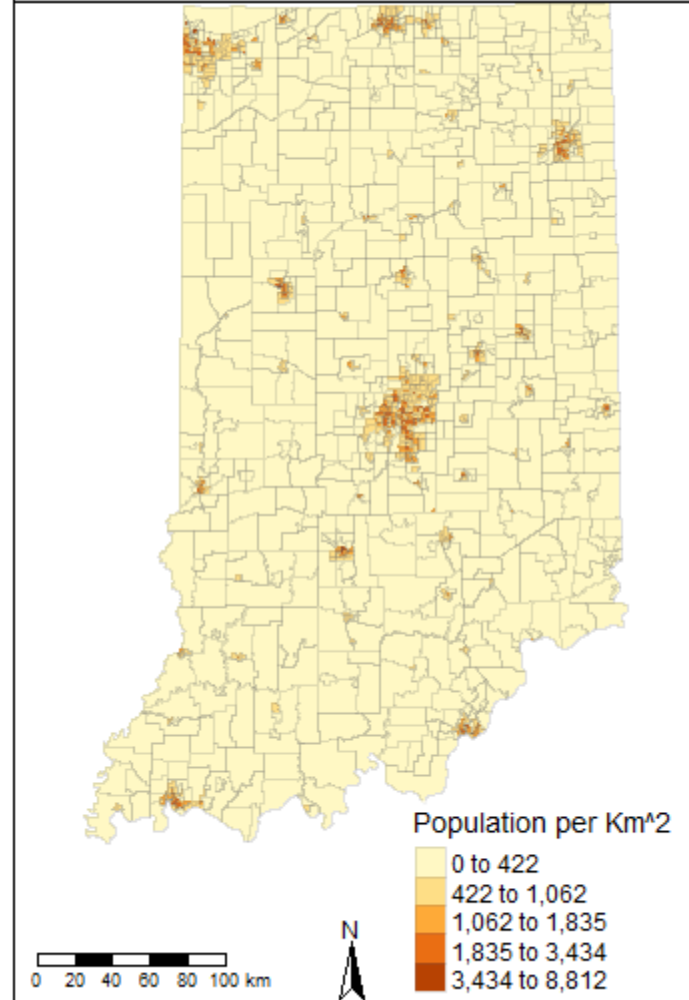
The following clustering methods were used Kmeans, Jenks, Hiearchical, Bayesian, Quantile, and Fisher

```
138  png('./figures/builtInClustering1.png', width=713, height=545)
139  tmap_arrange(k,j, nrow=1)
140  dev.off()
141  png('./figures/builtInClustering2.png', width=713, height=545)
142  tmap_arrange(h,b, nrow=1)
143  dev.off()
144  png('./figures/builtInClustering3.png', width=713, height=545)
145  tmap_arrange(q,f, nrow=1)
146  dev.off()
```
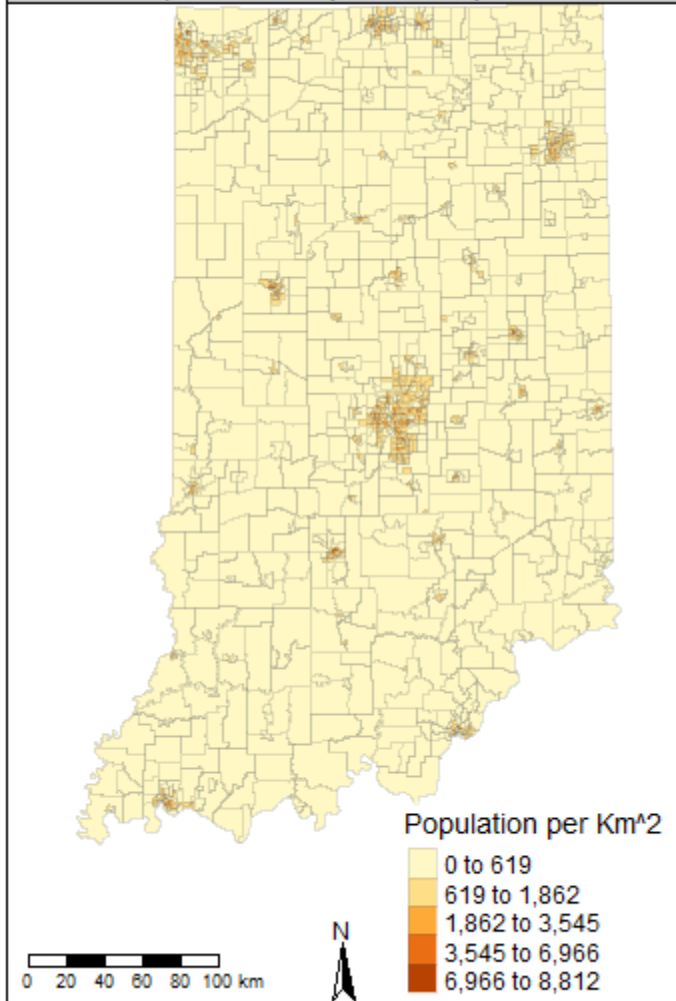
PURDUE
UNIVERSITY

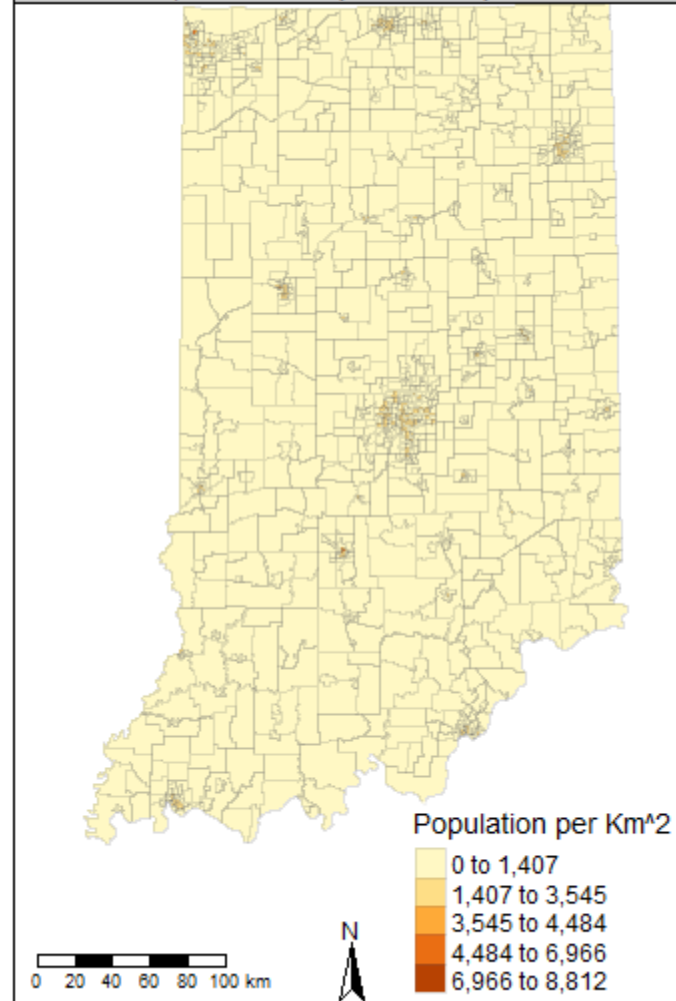# Developing a Kmeans Clustering Algorithm

```r
163  # the data value to be provided should be the item you want classes for
164  my_kmeans <- function(data, k) {
165    data <- as.numeric(data) # sanitize the data in the event that it has units
166    # first sort the data and get only the unique values
167    # you probably don't want the unique values
168    # unique values only will mess with the variance probably.
169    sorted_vals <- sort(data, decreasing=F)
170    # choose k elements at random, this will only happen once to initialize
171    centers <- sample(unique(sorted_vals), k) # use unique to guarantee no overlap
172    F_val.old <- -9999 # initializing the old F value to determine when to terminate the loop
173    F.all <- vector() # put in place to print an iterations graphic
174    rank.old <- 0 # initializing a container for previous iteration rankings.
175
176    # In my testing, 100 iterations seemed quick and suffient
177    for (j in seq(1, 100)) {
178
179
180      # need the distances for each k class
181      # line below will get the distance from centers and each value
182      center_dists <- t(sapply(sorted_vals, function(x) abs(centers-x)))
183
184      # find all values for each center that are less than evertyhing else
185      rank <- matrix(apply(center_dists, 1, which.min),ncol=1)
186      rank <- cbind(sorted_vals, rank)
187      colnames(rank) <- c('val', 'class')
188
189      # now, each row has a class from 1 to k assigned to it.
190      # group by the class and calculate variances
191      inClass.var = list()
```

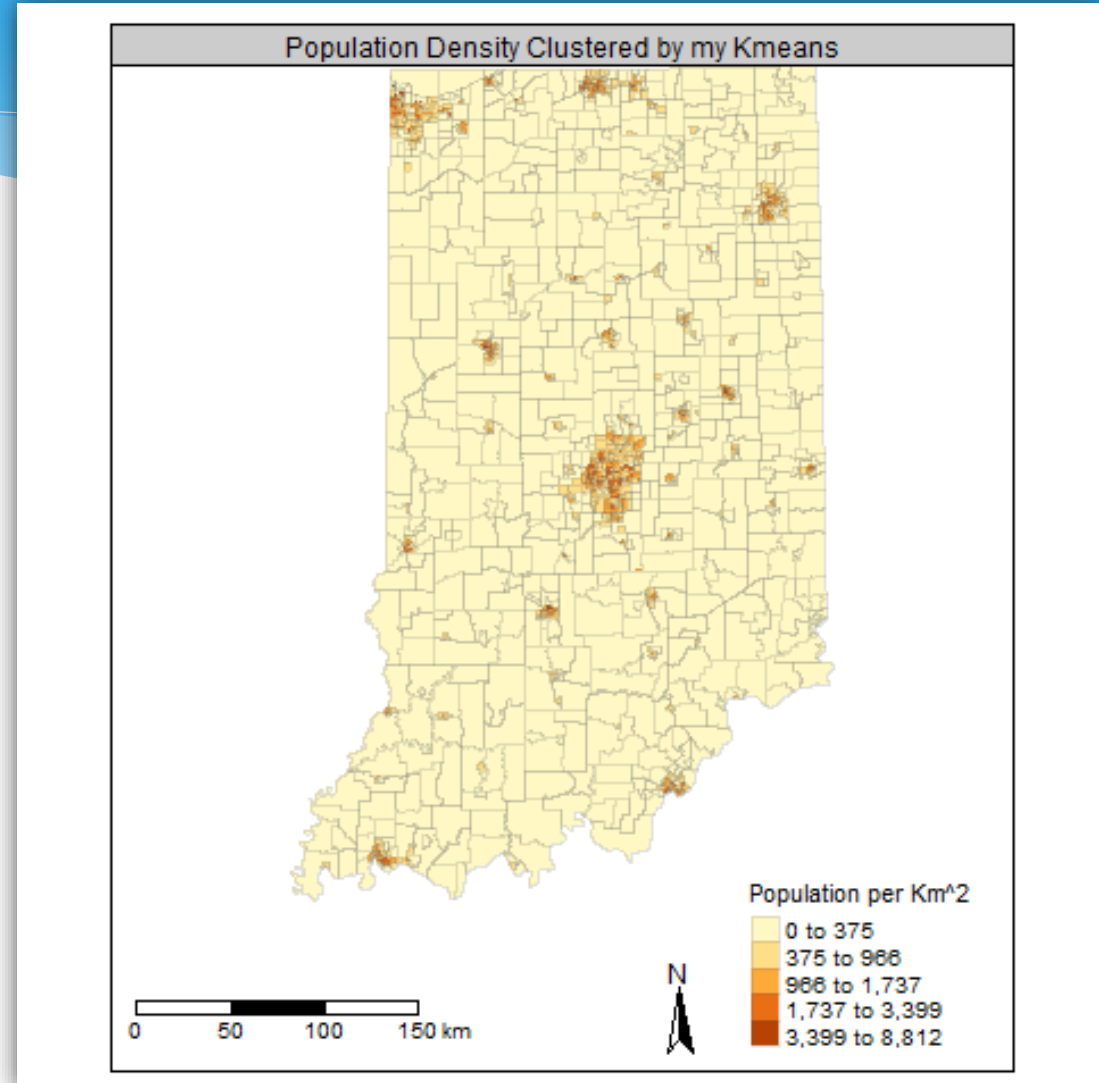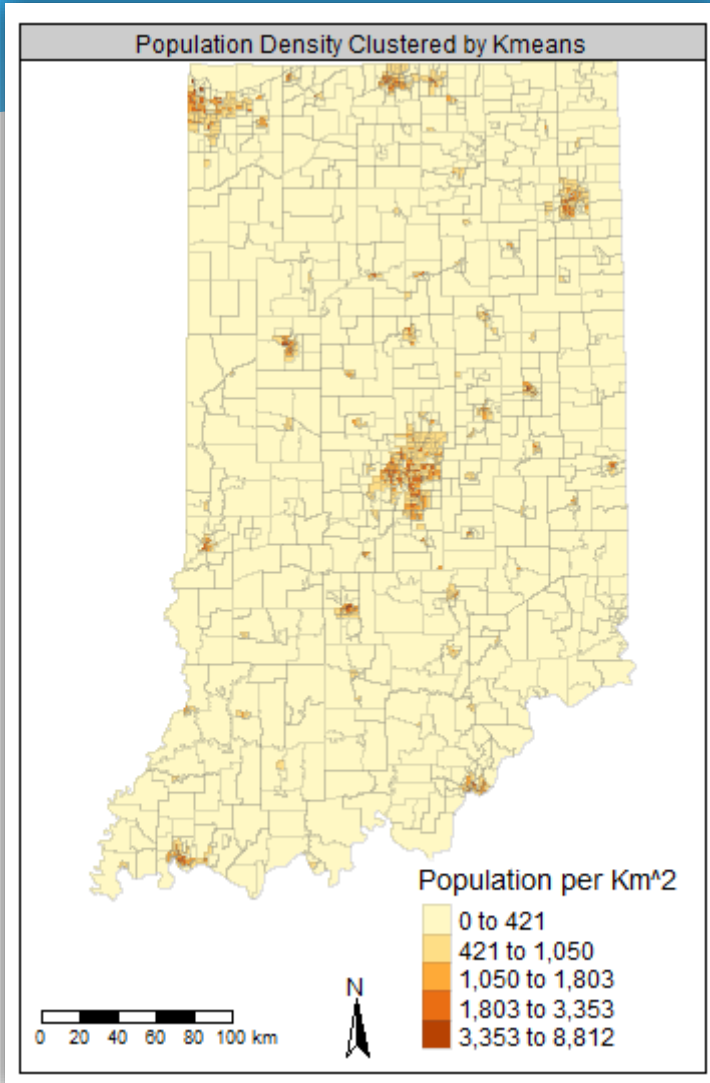# Developing a Kmeans Clustering Algorithm

```
193    for (i in 1:length(centers)) {
194      # grab all class values of a single type
195      inClass <- data.frame(rank) %>% dplyr::filter(class == i)
196      # find the variance in class
197      inClass.var[[i]] = var(inClass$val)
198
199      # adjust the class center
200      centers[i] = mean(inClass$val)
201    }
202
203    inClass.var <- sum(unlist(inClass.var))
204    betweenClass.var = var(centers)
205    F_val <-  betweenClass.var/inClass.var
206
207    F.all <- append(F.all, F_val)
208
209    # Save the largest version of F and the rankings.
210    # i.e. the jenk's criterion
211    if (F_val > F_val.old) {
212      iteration_max <- j
213      F_val.old <- F_val
214      rank.old <- rank
215    }
216
217  }
```

```
218      # rank.old should have the maximal F rankings
219      # break into a 2 column set with min and max values for each
220      classBreaks <- data.frame(rank.old) %>% group_by(class) %>%
221        summarize(min_val=min(val), max_val=max(val))
222
223      classBreaks <- as.vector(list(sort(c(classBreaks$min_val, max(classBreaks$max_val)))))[[1]])
224
225      png('./figures/F_graph.png')
226      plot(F.all, xlab='Iterations', ylab='F Value', main='Variation of F for each Iteration',
227          pch=16, col='deepskyblue3')
228      points(iteration_max, F_val.old, pch=16, col='red')
229      dev.off()
230
231      return(classBreaks)
232  }
```

PURDUE
U N I V E R S I T Y

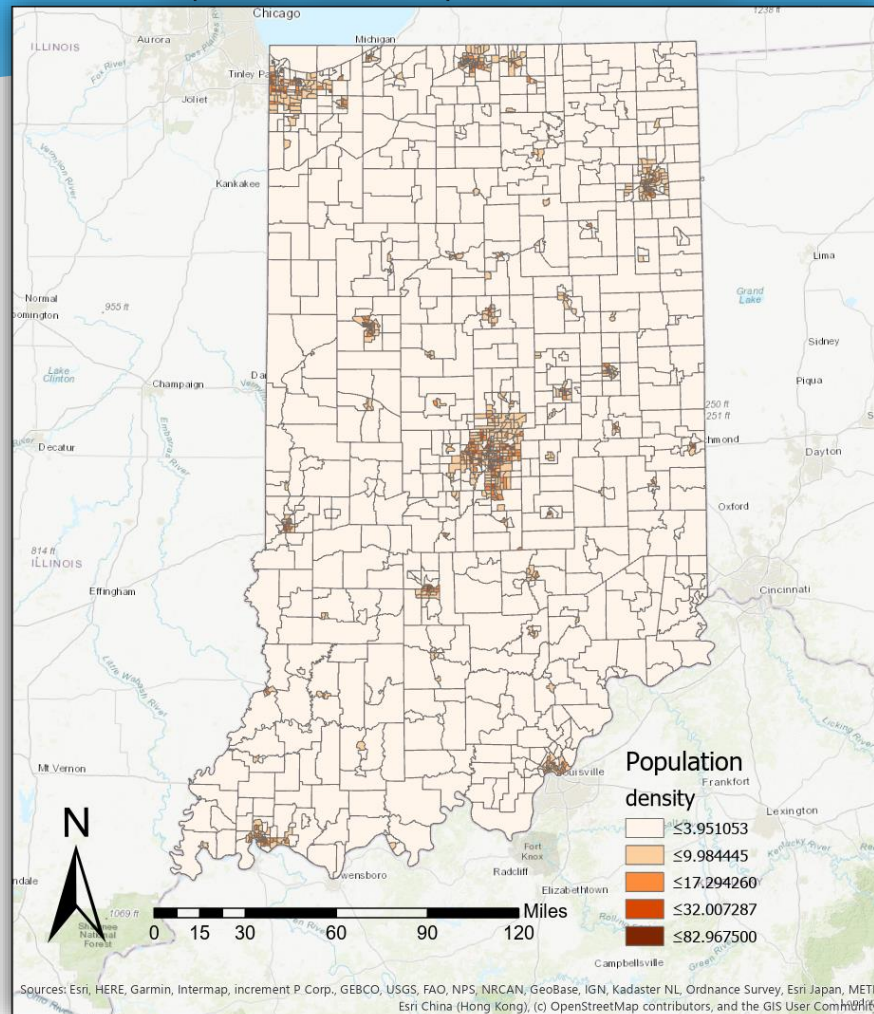# Developing a Kmeans Clustering Algorithm

```
239 # Tm will automatically fill in the gaps between the breaks
240 # since my algorithm will not ensure that each segment butts up against each other.
241 png('./figures/myKmeans.png')
242 tm_shape(tract_proj) +
243   tm_polygons('density', style='fixed', breaks=my_breaks,
244             title='Population per Km^2', border.alpha=.2 ) +
245   tm_compass(position=c('right', 'bottom'), text.size=1) +
246   tm_scale_bar(position=c('left', 'bottom'), text.size=.8) +
247   tm_layout(frame=T, panel.show=T, panel.labels='Population Density Clustered by my Kmeans',
248           legend.bg.color='white', legend.text.size=.8,
249           legend.bg.alpha=.6, inner.margins = c(.17, .2, 0, .2))
250 dev.off()
```

PURDUE
UNIVERSITY

# Developing a Kmeans Clustering Algorithm

# My Kmeans implementation vs ArcGIS

Population Clustered by Jenk's Natural Breaks