

An analysis of in-country article references and predictive modeling of local ratios based on country neighbors

Victoria Dollar, Ben Galin, Tyler Karhoff, and Gregory Smith

Emporia State University

January 13, 2023

Introduction and Background

Citations are an important tool for recording the use of public research. They allow recognition to authors for their methodology and contribution to their field. Citation formatting grants scholars the means to trace sources shared throughout the literature in a more manageable way. With the invention of the world wide web, published articles have become more accessible at greater distances. As a result, global inquiries examine the geolocation of the author in reference to subject matter, country of origin, number of citations used, the number of times the article was cited, etc.

An egocentric network centralizes the relationship between an entity of interest and its subnetwork. [9] For a particular publication, a local ratio can be constructed comparing the countries of the author to the countries of the citations' authors used in that article. A higher ratio would implicate that more citations utilized, originated from its own country whereas a lower ratio would reflect the usage of articles stemming from other countries. This paper will focus on comparing the ratios of various countries of statistics journals and estimating a country's ratio based upon the ratios of its neighbors. In the end, the hope is to gain insight on the relationship of networks between countries in the world of academia.

One factor that limits the use of a journal in research is accessibility. Since the birth of the internet, the number of journal publications have increased. Brazil (18.6%), the United Kingdom (10.7%), and the United States (6.4%) were reported to have the highest input in

open-access citations. [3] Open-access journals have contributed to this increase by allowing three-times more articles to be published, [5] and as a result, studies have been done to determine their impact. Antelman has reported that these journals have a much greater, yet complex impact on research. [2] It was found that the number of citations did not increase, but open-access articles were more favored. Similarly, Davis et al. notes that open-access allows for an increase in the number of readers but not necessarily the number of citations. [5] Although data will not be available to determine if open-access articles affect the local ratio, these articles could influence our results.

To date, no study has been published to determine the correlation between a country and its neighbors, in terms of academic references. However, one study was found that is similar to the study we performed. Gouel et al. examined the prevalence ratio, comparing the IP geolocations of readers to the IP geolocation of the author of online published articles. [6] The results showed that 98% of the webpage views derived from the country of origin. The mean prevalence ratio was 0.9 meaning that most of the articles did not network outside of their country. The minimum ratio was from Kosovo (0.18), and an unexpectedly low ratio of 0.6 was recorded for the United Kingdom. There were four additional studies that provided some insight into the number of citations used in articles in ocean acidification and fuzzy research studies by country.

Sahoo and Pandley tracked and mapped 100 published articles on ocean acidification. [7] The United States (57 citations), Germany (26 citations), and Australia (24 citations) were the top three countries that had their articles cited the most often. Alfaro-Garcia et al. published a study determining the total citation counts in fuzzy research articles for countries around the world, and the top three countries were China and Taiwan (141,298 citations), the United States (37,601 citations), and Iran (25,291 citations). [1]

Based upon this research, we can observe that these top countries are not located near one another. In effect, we can hypothesize that there will be no spatial correlation between a country's local ratio and the estimations produced from neighboring countries. Reasoning that the size, GDP, and commitment to education advancement will make a much larger impact with the number of publications and contribution to research over a country's geolocation.

Methods

Data collection

We started with an initial seed of the 100 most recent articles as of November 19, 2022, matching the query term `statistics` in Crossref, a DOI registration agency website. Identifying each article by its DOI, we then used Crossref’s public API to retrieve metadata about the article. Specifically, we retrieved metadata about the authors and their affiliations, and about the reference DOIs used by each article. We then repeated the process for each discovered reference DOI, running enough iterations until more than 500,000 articles were fetched.

For about 70% of articles, at least one of the authors included an affiliation name, often the name and address of a university. We attempted to find the country associated with each affiliation by parsing this value and matching it against a list of country names. We were able to associate about 8% of articles with a country in such a manner.

Detailed documentation for this data collection process is located in the Appendix under Listing 2.

Local ratio computation

For each article associated with a country, we looked at its references and their associated countries. We then computed the *article local ratio* as the number of references from the same country divided by the total number of references. We disregarded references that we were unable to associate with a country.

As a concrete example, consider this very article. Of the references listed below, seven have DOIs. Of these seven, only two — [7] and [4] — have affiliation metadata listed in the Crossref database, which we can parse and find to be India and Norway, respectively. Thus, this article’s local ratio is 0/2.

We then computed the *country local ratio* by taking the average of the article local ratios from this country. We are able to see these results in Figure 6 which was created using the R library `tmap`. [8]

Selecting neighbors

When selecting what constitutes a ‘neighbor’ for specific sovereign countries, we looked at three distinct methods. These methods were able to be performed using the R library `spdep`. [4] The first and perhaps most intuitive method is to say that neighbors are contiguous, or connected by a border. Doing so yields the result found in Figure 1.

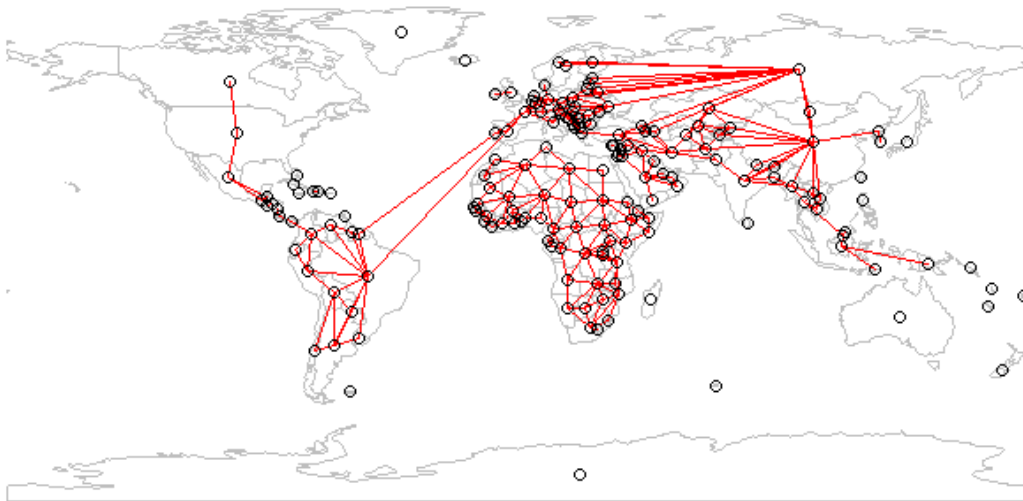


Figure 1: Boarder Neighbors (Red Line connects all Neighbors that Share a Border)

Notice that this is fairly accurate; however, it does have some drawbacks when it comes to analysis. For instance, countries with multiple sections cannot be well-matched. Notice that France is connected to Brazil and Suriname due to French Guiana being a part of France. This method's second problem is that island nations are all neighbor-less, making analysis with these countries impossible with this method. Thus, we will not be using this method for analysis, as it is too restrictive.

The second method we considered is using the centroids of all nations and a threshold maximal distance between them. There are many possible choices for a threshold. Two of them are shown in Figures 2 and 3.

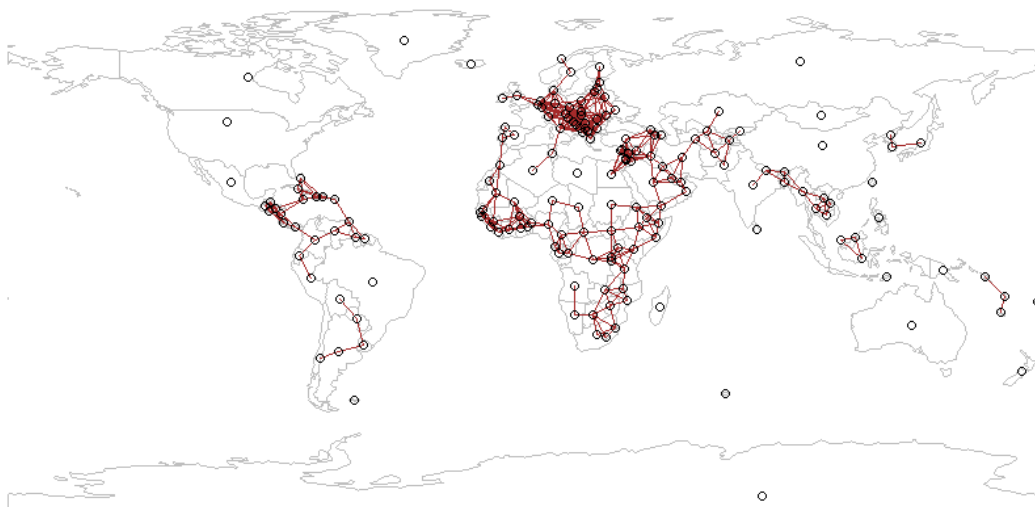


Figure 2: Centroid Neighbors (Brown Line connects all Neighbors that have close geographical centroids)

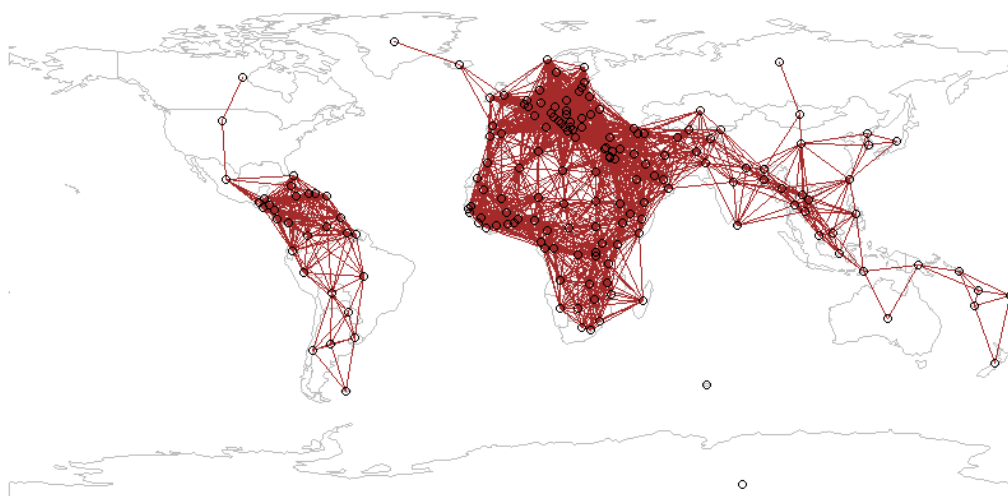


Figure 3: Centroid Neighbors (Brown Line connects all Neighbors that have further geographical centroids)

Both of these choices can cause problems with our data set. The first choice doesn't allow for much distance between neighbors and picks what countries are neighbors conservatively, while the second is much more liberal. We believe that the distinction in selecting neighbors with this method is also difficult to test properly, so we also have chosen not to use it. However, for a secondary study, it would be interesting to compare the results from this method and our last and chosen method.

The third method we considered is to use centroids as before, but instead of picking a maximal distance, we looked at a predefined number of nearest neighbors. As with the distance threshold approach, we experimented with different values for the number of nearest neighbors. Figures 4 and 5 show the results with three and five nearest neighbors.

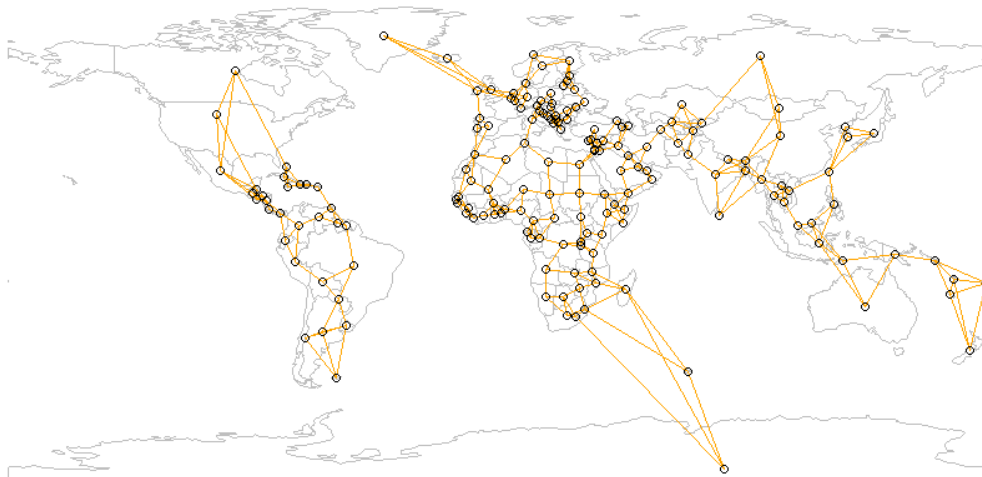


Figure 4: Centroid Neighbors (Orange Line connects three closest Neighbors)

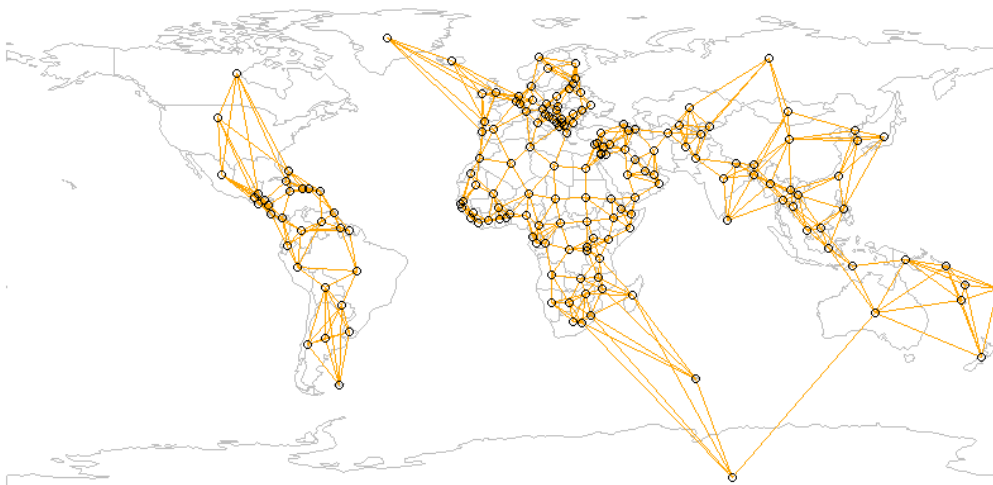


Figure 5: Centroid Neighbors (Orange Line connects five closest Neighbors)

The final method provides a pleasant result. We chose this method for our study for several reasons. First, unlike the first method, island nations can have neighbors for analysis. Second, varying the specific number of neighbors allows us to easily conduct several analyses: using a conservative choice for the number of neighbors and a more liberal choice for the

number of neighbors when trying to create our predictive model. Our hope was to see which neighbor count provided the best prediction for local ratios. If a neighboring country was missing its local ratio, then we selected the next neighbor. We found this to be the best way to pass over unknown countries without them vastly altering our results.

Spatial autocorrelation

Once we had a method of selecting neighbors, our goal was to see if there is spatial autocorrelation present in the country's local ratios. We computed the Moran I statistics [4] with three, five, and ten neighbors. We then adjusted the weighting function of the neighbors to an Inverse Distance Weight (IDW) and repeated the test.

Further, we tried to predict countries' local ratios by examining the local ratios of the countries nearest to them. Again, we used both a simple moving average and an IDW-based moving average for every sampled country. For example, the United States' closest three neighbors with data points are Canada, Mexico, and Cuba. Thus, when calculating the estimated ratio for the United States, we took the mean and the weighted averages of these three results. We then computed the Pearson correlation coefficient, between our sample local ratios and the estimated ratios. We did this for several neighbor counts to see if any provided a correlation. The details are in our results and analysis, and the code implementation is documented in Listing 1.

Listing 1: R Implementation of IDW

```

1 library(sf)
2 library(magrittr)
3
4
5 idw = function(Countries_df) {
6   World_centroids <- suppressWarnings(
7     st_centroid(Countries_df) %>%st_coordinates()
8   )
9
10  corr_list <- list()
11  for (m in 1 : 20) {
12    # Find the m nearest neighbors based on the centroids of each
13    # country.
14    World.knb <- knn2nb(knearneigh(World_centroids, k = m))
15
16    # 'country_sma' will hold the estimated local ratios, based on
17    # neighbors.

```

```

16   country_sma <- list()
17
18   # Iterate over countries.
19   for (i in 1 : nrow(Countries_df)) {
20     neighbor_list <- World.knb[[i]]
21
22     # Compute the first-power IDW of the current country.
23     # The formula is  $(\sum_j r_j/d_j)/(\sum_j 1/d_j)$ , where
24     #    $r_j$  is the local ratio of neighbor  $j$ ; and
25     #    $d_j$  is the centroid distance between neighbor  $j$  and our
country i.
26     numerator = 0
27     denominator = 0
28     for (j in 1 : m) {
29       distance <- as.numeric(st_distance(
30         suppressWarnings(st_centroid(Countries_df[i,])),
31         suppressWarnings(st_centroid(Countries_df[neighbor_list[[j
32         ]],])))
33       local_ratio <- Countries_df[neighbor_list[[j]],]$local_ratio
34       numerator <- numerator + local_ratio/distance
35       denominator <- denominator + 1/distance
36     }
37     country_sma <- append(country_sma, numerator/denominator)
38   }
39
40   # Compute the Pearson correlation between the actual local_ratios
and the
41   # estimated ones.
42   country_sma <- as.numeric(country_sma)
43   knb_final_df <- data_frame(Countries_df$sovereight, Countries_df$
local_ratio, country_sma)
44   colnames(knb_final_df) <- c('sovereight', 'local_ratio', 'country_
sma')
45   corr_list <- append(corr_list, cor(knb_final_df$local_ratio, knb_
final_df$country_sma))
46 }
47 return(corr_list)
48 }

```


Results & Analysis

The results located in Table 1 compare the sample local ratios with their estimations based on a simple moving average of the neighbors' local ratios. We are showing the results for three-, five-, and ten-closest neighbors, along with visualizations for each in figures 7, 8, and 9.

Country	Sample Local Rat.	3 Neigh Est.	5 Neigh Est.	10 Neigh Est.	Country	Sample Local Rat.	3 Neigh Est.	5 Neigh Est.	10 Neigh Est.
Algeria	0.0000	0.1791	0.1688	0.1036	Luxembourg	0.0000	0.1502	0.1588	0.1735
Argentina	0.1855	0.0738	0.0443	0.0681	Malaysia	0.0286	0.0787	0.0472	0.1066
Australia	0.1893	0.0095	0.0331	0.0991	Luxembourg	0.0000	0.1502	0.1588	0.1735
Austria	0.1084	0.1396	0.1518	0.1140	Malaysia	0.0286	0.0787	0.0472	0.1066
Bangladesh	0.3333	0.1217	0.1031	0.1017	Mexico	0.0189	0.3710	0.2393	0.1769
Belgium	0.1288	0.1073	0.1330	0.1606	Morocco	0.3667	0.0569	0.0654	0.0805
Brazil	0.1007	0.0000	0.0612	0.0747	Mozambique	0.0000	0.1596	0.1237	0.1114
Canada	0.1658	0.3773	0.2430	0.1531	Myanmar	0.0000	0.1898	0.1698	0.1350
Chile	0.1206	0.0954	0.0572	0.0746	Netherlands	0.1715	0.1062	0.1245	0.1564
China	0.1505	0.1111	0.1412	0.1088	Nigeria	0.0000	0.0000	0.0854	0.1318
Colombia	0.0000	0.0278	0.1082	0.1622	Norway	0.1859	0.1684	0.1042	0.1516
Costa Rica	0.0000	0.1468	0.0919	0.1622	Pakistan	0.3356	0.0934	0.1561	0.1089
Croatia	0.0000	0.0842	0.0806	0.1219	Panama	0.0833	0.1190	0.0752	0.1539
Cuba	0.3571	0.0278	0.0205	0.1310	Philippines	0.0000	0.0551	0.0803	0.1141
Cyprus	0.0000	0.1094	0.0870	0.0885	Poland	0.1282	0.1599	0.1176	0.0835
Czech Republic	0.3772	0.0500	0.0762	0.1025	Portugal	0.0677	0.2087	0.1876	0.1389
Denmark	0.1535	0.1204	0.1735	0.1590	Qatar	0.1667	0.0770	0.0462	0.1124
Egypt	0.0601	0.1094	0.0657	0.0675	Romania	0.0000	0.0878	0.0783	0.0919
Estonia	0.0156	0.1101	0.0960	0.1146	Russia	0.1436	0.1880	0.1799	0.1577
Finland	0.2020	0.0551	0.0959	0.1113	Saudi Arabia	0.0798	0.0556	0.0662	0.0706
France	0.1563	0.0730	0.1699	0.1300	Senegal	0.0000	0.1222	0.0869	0.1043
Georgia	0.2288	0.0356	0.0213	0.0666	Slovakia	0.0000	0.2027	0.1299	0.1098
Germany	0.1898	0.1073	0.1167	0.1282	Slovenia	0.0417	0.1619	0.1477	0.1207
Greece	0.1607	0.0342	0.0288	0.0887	South Africa	0.0690	0.1366	0.0940	0.0880
Guinea	0.0000	0.1222	0.0869	0.1043	South Korea	0.0801	0.1544	0.0926	0.1061
Hungary	0.1026	0.0139	0.1094	0.1156	Spain	0.1030	0.1969	0.1740	0.1353
Iceland	0.1667	0.1561	0.1537	0.1279	Sweden	0.1498	0.1805	0.1114	0.1424
India	0.1291	0.2230	0.1810	0.1453	Switzerland	0.1504	0.1062	0.1154	0.1321
Indonesia	0.0000	0.0095	0.0803	0.1155	Taiwan	0.1367	0.0267	0.0813	0.1005
Iran	0.1511	0.0556	0.0950	0.1139	Thailand	0.2361	0.1111	0.0997	0.0858
Iraq	0.0000	0.0547	0.0816	0.1061	Turkey	0.1067	0.0547	0.0657	0.0718
Ireland	0.0484	0.1830	0.1576	0.1362	Ukraine	0.0000	0.0356	0.0245	0.0743
Israel	0.1642	0.0547	0.0449	0.0800	United Kingdom	0.2637	0.1162	0.1010	0.1290
Italy	0.1503	0.0500	0.1355	0.1131	United Arab Emirates	0.0000	0.1325	0.1466	0.1290
Japan	0.1759	0.0723	0.0735	0.0965	United States	0.7558	0.1806	0.1250	0.1013
Jordan	0.0000	0.1094	0.0657	0.0970	Uruguay	0.0000	0.1356	0.0814	0.0847
Kenya	0.4097	0.0466	0.0613	0.0704	Venezuela	0.0000	0.0278	0.1082	0.0866
Latvia	0.0000	0.1153	0.0692	0.1084	Vietnam	0.0000	0.0883	0.0803	0.1095
Lebanon	0.0000	0.1094	0.0870	0.0964					

Table 1: Sample Ratios and Estimated Local Ratios for 3, 5, and 10 Neighbors

The estimates tend to be more accurate when the local ratio is around 0.10-0.20. Exceptionally high or exceptionally low local ratios are estimated poorly by their neighbors. This is an example of a “regression to the mean”, as the neighbors’ local ratios are less extreme.

Looking at the top ten countries based on local ratios, we didn’t notice any obvious trends. These countries — United States, Kenya, Czech Republic, Morocco, Cuba, Pakistan, Bangladesh, United Kingdom, Thailand, and Georgia — include both large and small countries, countries from different continents, high GDP and low GDP countries, etc.

Using the Moran I test, we found no autocorrelation in the local ratios, irrespective of the number of neighbors and the average weighting function used. See Table 2.

Neighbors	Weight	p -value	Moran I statistic
3	simple average	0.673	-0.047723953
3	IDW	0.7221	-0.08322051
5	simple average	0.5844	-0.02575441
5	IDW	0.6503	-0.05582828
10	simple average	0.6703	-0.030992166
10	IDW	0.5325	-0.02120910

Table 2: Moran I test results

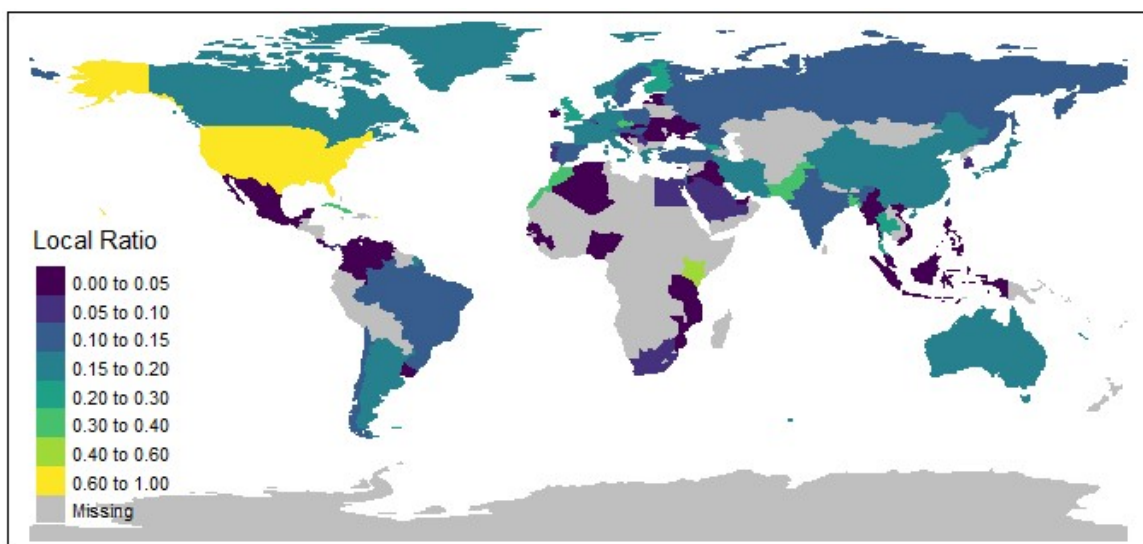


Figure 6: Local Ratio by Country (Breaks are not perfectly divided. This is present in all future maps)

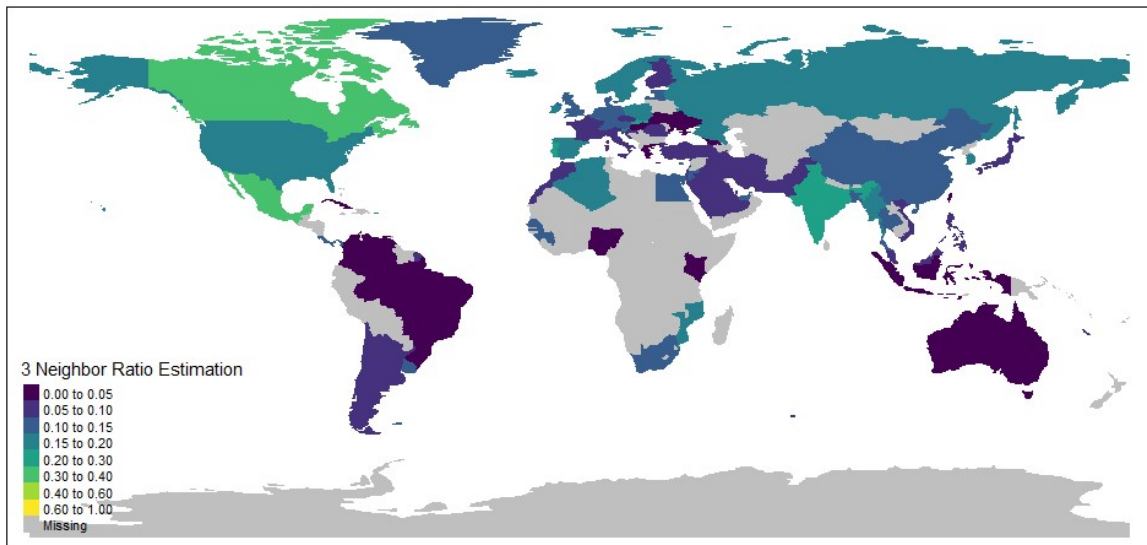


Figure 7: Estimating Local Ratio Based on 3 Closest Neighbors

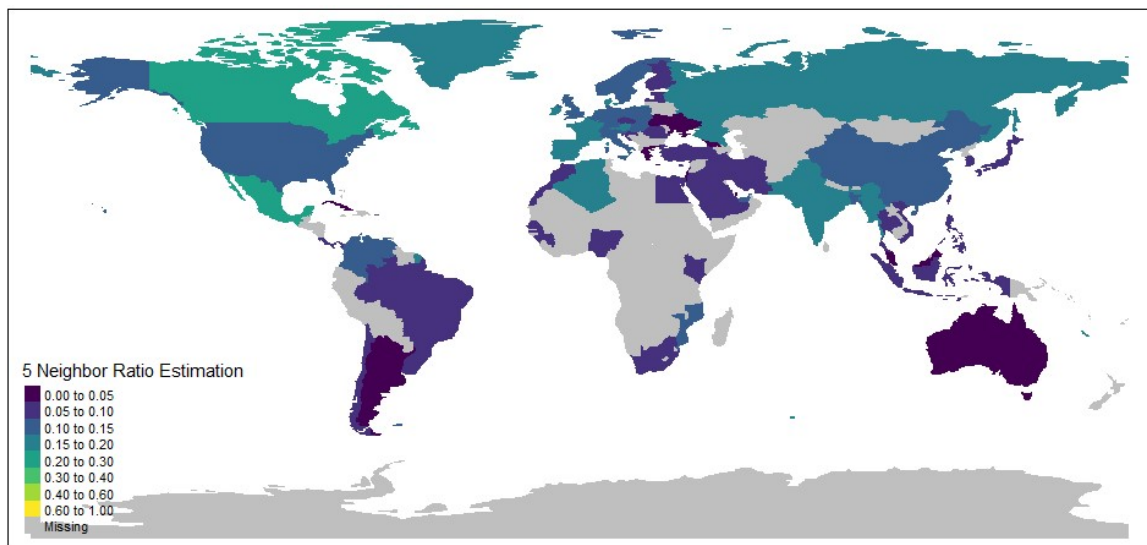


Figure 8: Estimating Local Ratio Based on 5 Closest Neighbors

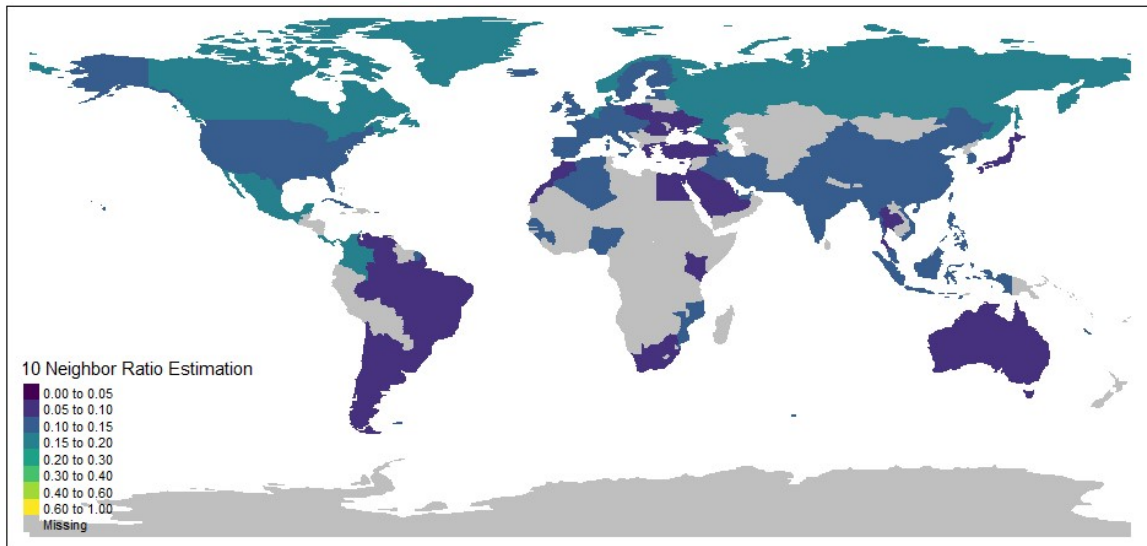


Figure 9: Estimating Local Ratio Based on 10 Closest Neighbors

The thematic maps [8] in Figures 7, 8, and 9 visualize the estimated local ratios for three-, five-, and ten-neighbors, respectively. The Pearson correlation coefficients are -0.08678 for three-, -0.07008 for five-, and -0.14110 for ten neighbors. When peering into more neighbors the Pearson coefficient stays negative with little correlation. Similarly, the Moran I test confirms no autocorrelation is present. This shows that the sample local ratio cannot be reliably estimated from the local ratios of neighboring countries.

Discussion & Conclusion

As we hypothesized, there is no correlation between a country's local ratio and the estimation of the neighboring countries. We can conclude that using k -neighbors, is not the best method of estimation per the results of the Moran I test. In respect to this, there are many contributing factors that could influence the lack of correlation, such as, the number of universities and other research facilities located in a country, the country's GDP and population, and geolocation in reference to other research-focused groups.

In addition, the process of matching an article to a country is imprecise. The affiliation data stored with each article is open-form, and authors often enter non-standard address information or location. Since we processed a large number of records, we had to automate this parsing and matching logic. In the case of United States-based articles, we were aided by our knowledge of common forms of address writing (for example, addresses ending with a state two-letter abbreviation and a 5- or 9-digit zip code). Presumably, we were thus able to

detect a larger proportion of United States articles. This, in effect, contributed to a higher computed local ratio in the United States as compared to the true local ratio. Conversely, for countries with a large fraction of United States citations, we computed a lower local ratio than we would have had, had we not been able to parse United States addresses. A more sophisticated language parsing scheme, such as from natural language processing, would undoubtedly improve this result, which is an avenue of potential future study.

Another main contributor to the data quality is the sparse nature in which reference DOIs are present for an article given the different citation forms (APA, MLA, Chicago, AMS, etc.). Consequently, many articles do not have any importance in regard to analysis besides being a reference to source article. Ultimately, both of the above factors require retrieval of mass articles to gather a sufficient sample of articles that can be assigned to source countries. For a more representative sample, magnitudes greater than half a million articles should be amassed for ratio analysis.

References

- [1] ALFARO-GARCIA, V., AND MERIGO, J. A citation analysis of fuzzy research by universities and countries. *Journal of Intelligent & Fuzzy Systems* 38 (2020).
- [2] ANTELMAN, K. Do open-access articles have a greater research impact? *College & Research Libraries* 65, 5 (2004).
- [3] BHARDWAJ, R. India's contribution to open access movement. *Journal of Knowledge & Communication Management* 5, 2 (2015), 107–126.
- [4] BIVAND, R. R packages for analyzing spatial data: A comparative case study with areal data. *Geographical Analysis* 54, 3 (2022), 488–518.
- [5] DAVIS, P., LEWENSTEIN, B., AND ET AL. Open access publishing, article downloads, and citations: Randomised controlled trial. *BMJ* (2008).
- [6] GOUEL, M., VERMEULEN, K., AND ET AL. Longitudinal study of an ip geolocation database. *arXiv* (2021).
- [7] SAHOO, S., AND PANDEY, S. Characteristics and inter-citation network of 100 most influential studies on ocean acidification: A bibliometric analysis. *Science & Technology Libraries* 41, 1 (2022), 56–72.

- [8] TENNEKES, M. tmap: Thematic maps in R. *Journal of Statistical Software* 84, 6 (2018), 1–39.
- [9] WENG, T., LI, Z., AND ZHANG, J. Egocentric visual analysis of dynamic citation network. *Journal of Visualization* 25 (2022), 1343–1360.

Appendices

Python Code

```

1 import re
2 import os
3 import pandas as pd
4 import numpy as np
5 import time
6 from concurrent.futures import ThreadPoolExecutor, as_completed
7 from habanero import Crossref
8 import urllib.parse
9 import math
10 from typing import Set, Dict, List, Tuple
11
12
13 # ARTICLE RETRIEVAL #
14
15 def parse_dois(article: dict) -> Set[str]:
16     return set([urllib.parse.unquote(ref.get('DOI').strip()) for ref in
17                 article.get('reference', []) if ref.get('DOI')])
18
19 def fetch_batch_to_df_threaded(dois: List[str], max_articles=100,
20                                max_iteration=3,
21                                runner_sleep=0.1, max_workers=20) -> pd.DataFrame:
22     articles_df = pd.DataFrame()
23     start_time = time.time()
24     iteration = 0
25     while (articles_df.shape[0] < max_articles) and (len(dois) > 0):
26         iteration += 1
27         # boolean that determines if we get next 'dois'
28         # on final iteration is 'False' so len(dois) == 0 and loop
29         # breaks
30         get_next = iteration < max_iteration
31         dois, articles_df = fetch_threaded(dois, articles_df,
32                                             start_time, get_next, runner_sleep, max_workers, max_articles)
33         print('== end of iteration #s (%d articles) (%d batches next
34               iteration) (%.2fs) =='
35               % (iteration, articles_df.shape[0], math.ceil(len(dois)
36                    /1000), time.time() - start_time))
37     print('Total articles:', articles_df.shape[0])
38     print('Final iteration:', iteration)
39     print('Final time: %.2fs' % (time.time() - start_time))

```

```

35     print()
36     return articles_df
37
38
39 def fetch_threaded(dois: List[str], articles_df: pd.DataFrame,
40     start_time: time.time, get_next: bool,
41     runner_sleep: float, max_workers: int, max_articles: int) ->
42     Tuple[List[str], pd.DataFrame]:
43     next_dois = set()
44     batch = 0
45     while dois:
46         # we get some DOIs
47         dois_temp = dois[:1000]
48         # request data for each DOI
49         futures = runner(dois_temp, runner_sleep, max_workers)
50
51         # once we're received all our DOI requests we can proceed
52         # we create the articles dictionary for these 1000 articles
53         # and add the corresponding reference DOIs to a list to be
54         # parsed in next iteration of caller
55         articles = {}
56         articles_togo = max_articles - articles_df.shape[0] - len(dois)
57         queried_dois = set(articles_df.index).union(set(dois))
58         for future in as_completed(futures):
59             doi_temp, dict_temp = future.result()
60             articles[doi_temp] = dict_temp
61             # if we're not parsing the next ancestor level, let's save
62             some time;
63             # also, once we have enough total articles, let's not look
64             for more ancestors
65             if (get_next) and (len(next_dois) < articles_togo):
66                 next_dois = next_dois.union(parse_dois(dict_temp).
67                 difference(queried_dois))
68
69         batch += 1
70         # create corresponding DataFrame for articles and add to
71         ongoing DataFrame
72         articles_df_temp = articles_json_to_df(articles)
73         articles_df = articles_df.append(articles_df_temp)
74         print('** Duplicates dropped: %d **' % articles_df[articles_df.
75         index.duplicated()].shape[0])
76         articles_df = articles_df[~articles_df.index.duplicated(keep='

```

```

first')]
    # articles_df.to_csv('backup.csv', index=True)
    print('** end of batch      #d (%d articles) (%d batches next
iteration) (%.2fs) **'
          % (batch, articles_df.shape[0], math.ceil(len(next_dois)
/1000), time.time() - start_time))
    # remove the DOIs we just queried
    dois = dois[1000:]

    return list(next_dois), articles_df

79 def runner(doi_list: List[str], runner_sleep: float, max_workers: int)
-> List:
80     futures = []
81     cr = Crossref(mailto = "your@email.ext")
82     with ThreadPoolExecutor(max_workers=max_workers) as executor:
83         for doi in doi_list:
84             futures.append(executor.submit(get_article_with_etiquette,
doi, cr, runner_sleep))
85             time.sleep(runner_sleep)
86     return futures

87
88
89 def get_article_with_etiquette(doi: str, cr: Crossref, runner_sleep:
float) -> Tuple[str, Dict]:
90     count = 0
91     while count < 100:
92         try:
93             r = cr.works(ids = doi)
94         except Exception as e:
95             error_substring = str(e).split(':')[1].strip() if len(str(e)
).split(':') > 1 else str(e)
96             if error_substring == 'Too Many Requests for url':
97                 extra_long = 300
98                 print('Sleeping for a while; %d seconds' % extra_long)
99                 time.sleep(extra_long)
100             elif error_substring == 'Not Found for url':
101                 break
102             elif error_substring in ['Max retries exceeded with url', '
Gateway Time-out for url']:
103                 print(error_substring, doi)

```

```

104         else:
105             print(e)
106             time.sleep(runner_sleep)
107             count += 1
108         else:
109             # bundle the article data with corresponding DOI, since '
as_completed'
110             # iterates based upon completion and not list order
111             return doi, r.get('message')
112     return doi, {}
113
114
115 # ARTICLE DATA TRANSFORMATION #
116
117 def get_key_keep_list() -> List:
118     return ['DOI', 'author', 'reference']
119
120
121 def handle_prop_key(prop_key_str: str, prop_val):
122     if prop_key_str == 'author':
123         prop_list = list()
124         # create a list, consisting of all
125         # affiliations related to article
126         for author in prop_val:
127             affn_temp = author['affiliation']
128             if len(affn_temp) > 0:
129                 prop_list.extend([x.get('name') for x in affn_temp if x
.get('name')])
130             if len(prop_list) > 0:
131                 return prop_list
132         else:
133             return None
134
135     elif prop_key_str == 'reference':
136         prop_list = list()
137         # add all reference DOIs to a list
138         for reference in prop_val:
139             doi_temp = reference.get('DOI')
140             if doi_temp:
141                 prop_list.append(doi_temp)
142         if len(prop_list) > 0:
143             return prop_list

```

```

144     else:
145         return None
146
147     elif prop_key_str == 'published':
148         date_list_nested = prop_val.get('date-parts')
149         # return date in YYYY-MM-DD string format
150         if date_list_nested:
151             prop_val = ['-'.join([str(y) for y in x]) for x in
date_list_nested]
152         else:
153             return date_list_listed
154
155     # 'DOI', 'subject', 'title', 'container-title', 'reference-count'
and non-null 'published' date
156     # fall through to here
157     if type(prop_val) == int:
158         return prop_val
159     elif len(prop_val) == 0:
160         return None
161     # if singleton list, just return the element
162     elif isinstance(prop_val, list) and len(prop_val) == 1:
163         return prop_val[0]
164     else:
165         return prop_val
166
167
168 def article_list_for_df(article_dict: Dict) -> List[dict]:
169     key_keep_list = get_key_keep_list()
170     article_list = list()
171
172     # iterate through the articles
173     for doi, prop_keys in article_dict.items():
174         # iterate through the keys of the article
175         if prop_keys:
176             temp_dict = dict()
177             for prop_key in prop_keys:
178                 if prop_key in key_keep_list:
179                     temp_dict[prop_key] = handle_prop_key(prop_key,
prop_keys[prop_key])
180             else:
181                 temp_dict = {'DOI': doi}
182             article_list.append(temp_dict)

```

```

183
184     return article_list
185
186
187 def articles_json_to_df(articles_json: json) -> pd.DataFrame:
188     # dataframe creation
189     articles_list = article_list_for_df(articles_json)
190     articles_df = pd.DataFrame(data = articles_list, columns =
191     get_key_keep_list())
192
193     # some basic cleaning
194     df_cols = articles_df.columns
195     articles_df.fillna(np.nan, inplace=True)
196     articles_df.set_index('DOI', inplace=True)
197     articles_df = articles_df[~articles_df.index.duplicated(keep='first
198     ')]
199
200     # basic variable creation and data type fixing
201     articles_df['reference_len'] = articles_df['reference'].apply(
202     lambda x: get_reference_len(x)).astype(int)
203     articles_df['published'] = pd.to_datetime(articles_df['published'],
204     errors = 'coerce')
205     articles_df['year'] = articles_df['published'].dt.year
206     articles_df['year'] = articles_df['year'].astype('Int64')
207     articles_df.drop(['published'], axis=1, inplace=True)
208
209     ## we want to know if we queried all references, so we need to keep
210     all articles, ##
211     ## even if they don't have author affiliations ##
212     # primary affiliation to country mapping
213     articles_df['country'] = articles_df['author'].apply(lambda x:
214     get_primary_author_country(x))
215
216     return articles_df
217
218 def get_primary_author_country(name_list: List[str]) -> str:
219     if type(name_list) == list:
220         return clean_affiliation_name(name_list[0])
221     return np.nan

```

```

219 def get_reference_len(reference: List) -> int:
220     try:
221         return len(reference)
222     except TypeError:
223         return 0
224
225
226 # ARTICLE PRIMARY AFFILIATION CLEANING #
227
228 def clean_affiliation_name(name: str) -> str:
229     found = find_country_name_at_end_of_affiliation(name)
230     if found:
231         return found
232     found = find_state_name_at_end_of_affiliation(name)
233     if found:
234         return found
235     return None
236
237
238 def get_countries() -> Set[str]:
239     with open('Countries.csv', 'r') as handle:
240         countries = set([line.strip() for line in handle])
241     return countries
242
243
244 def get_country_synonyms() -> Dict[str, str]:
245     country_synonyms = {}
246     with open('Country_synonyms.csv', 'r') as handle:
247         for f in handle.readlines():
248             k, v = f.split(',')
249             country_synonyms[k.strip().lower()] = v.strip()
250     return country_synonyms
251
252
253 def get_us_states() -> Set[str]:
254     us_states = set()
255     with open('US_states.csv', 'r') as handle:
256         for f in handle.readlines():
257             name, abbr = f.split(',')
258             us_states.update([name.strip(), abbr.strip()])
259     return us_states
260

```

```

261
262 def find_country_name_at_end_of_affiliation(name: str) -> str:
263     countries = get_countries()
264     country_synonyms = get_country_synonyms()
265     name = name.rstrip(' ,;-().')
266     country_regex = '|'.join(sorted(countries, key=len, reverse=True))
267     match = re.search('(' + country_regex + ')$', name, re.I)
268     if match:
269         return match.group()
270     country_synonyms_regex = '|'.join(sorted(country_synonyms.keys(),
271     key=len, reverse=True))
272     match = re.search('(' + country_synonyms_regex + ')$', name, re.I)
273     if match:
274         try:
275             return country_synonyms[match.group().lower()]
276         except KeyError:
277             print('Synonym not found:', match.group().lower())
278     return ''
279
280 def find_state_name_at_end_of_affiliation(name: str) -> str:
281     us_states = get_us_states()
282     name = name.rstrip(' ,;-().')
283     us_states_regex = '|'.join(us_states)
284     zip_code_regex = '[0-9]{5}(-[0-9]{4})?'
285     match = re.search('(' + us_states_regex + ')[, ]*' + zip_code_regex
286     + '$', name, re.I)
287     if match:
288         return 'United States'
289     return ''
290
291 # ARTICLE AGGREGATION #
292
293 def get_articles_from_csv(filename: str, dir='articles') -> pd.
294 DataFrame:
295     path = '/'.join([dir, filename])
296     articles_df = pd.read_csv(path, index_col=0)
297     if articles_df.empty:
298         return articles_df
299
300 # if needed, we need to rebuild the lists from string

```



```

representation
300     columns = articles_df.columns
301     if 'reference' in columns:
302         articles_df['reference'] = articles_df['reference'].apply(
lambda x: string_to_list(x))
303     if 'author' in columns:
304         articles_df['author'] = articles_df['author'].apply(lambda x:
string_to_first_list_element(x))
305     return articles_df
306
307
308 def string_to_list(str_of_list: str) -> list:
309     if type(str_of_list) == str:
310         return eval(str_of_list)
311     else:
312         return str_of_list
313
314
315 def string_to_first_list_element(str_of_list: str) -> str:
316     if type(str_of_list) == str:
317         return eval(str_of_list)[0]
318     else:
319         return str_of_list
320
321
322 def get_all_articles_from_csv() -> pd.DataFrame:
323     articles_list = os.listdir('articles')
324     articles_df = pd.DataFrame()
325     for index, articles in enumerate(articles_list):
326         print(index, articles)
327         articles_temp_df = get_articles_from_csv(articles)
328         articles_temp_df['country'] = articles_temp_df['country'].apply
(lambda x: x.title() if type(x) == str else x)
329         articles_df = articles_df.append(articles_temp_df)
330         articles_df = articles_df[~articles_df.index.duplicated(keep='
first')]
331     return articles_df
332
333
334 def get_cnty_counts_df(articles_input_df: pd.DataFrame) -> pd.DataFrame
:
335     # we create a new derivative DataFrame

```

```

336 articles_df = articles_input_df.drop(['author'], axis=1)
337
338 # default ratio value when no references present (refs == None)
339 articles_df['local_ratio'] = np.nan
340 # default count when no references present
341 articles_df['reference_used'] = 0
342
343 # we don't want to drop articles without a country,
344 # but we only want to consider source articles with a country
345 print(articles_df[~articles_df['country'].isna()].shape)
346 for index in articles_df[~articles_df['country'].isna()].index:
347     refs = articles_df.loc[index, 'reference']
348     # refs should either be non-empty list or np.nan
349     if type(refs) == list:
350         source_cnty = articles_df.loc[index, 'country']
351         refs_in_cnty, refs_total = get_cnty_counts_for_refs(refs,
articles_df, source_cnty)
352
353         # either: the article does not have all DOI ancestors in DF
, not a valid observation;
354         # the article has no references with a country, does not
need to be updated
355         if refs_total == 0:
356             continue
357
358         ratio_temp = refs_in_cnty / refs_total
359         articles_df.loc[index, 'local_ratio'] = ratio_temp
360         articles_df.loc[index, 'reference_used'] = refs_total
361
362     # although strictly not needed since np.na will be ignored for
grouped statistics by default,
363     # doing this will reduce any human error
364     articles_df.dropna(subset=['local_ratio'], inplace=True)
365     articles_df.drop(['reference'], axis=1, inplace=True)
366     # data type fix from float
367     articles_df['reference_used'] = articles_df['reference_used'].
astype(int)
368     return articles_df
369
370
371 def get_cnty_counts_for_refs(refs: List[str], articles_df: pd.DataFrame
, source_cnty: str) -> Tuple[int, int]:

```

```
372 refs_in_cnty = 0
373 refs_total = 0
374 for ref in refs:
375     ref = urllib.parse.unquote(ref.strip())
376     try:
377         cnty_temp = articles_df.loc[ref, 'country']
378         # reference not in DataFrame
379     except KeyError:
380         pass
381     # if successful
382     else:
383         # when entry is not np.nan
384         if cnty_temp == cnty_temp:
385             refs_total += 1
386         if cnty_temp == source_cnty:
387             refs_in_cnty += 1
388 return refs_in_cnty, refs_total
389
390
391 def create_grouped_dfs(articles_df: pd.DataFrame) -> pd.DataFrame:
392     articles_cnty_df = articles_df.groupby(['country']).mean()
393     return articles_cnty_df
```

Listing 2: Python Functions Developed For Data Collection, Cleaning, and Transformation

R Code

```
1 library(mapdata)
2 library(tidyverse)
3 library(sf)
4 library(tigris)
5 library(spData)
6 library(tmap)
7 library(cartogram)
8 library(isdas)
9 library(gridExtra)
10 library(plotly)
11 library(patchwork)
12 library(spdep)
13 library(GWmodel)
14 library(kableExtra)
15 library(spatialreg)
16 library(spgwr)
17
18
19 ##### Retrieve and Prepare Data #####
20
21 data('World')
22 articles <- read_csv('country_counts.csv')
23 colnames(articles) <- c('sovereight', 'ref', 'local_ratio', 'reference_
    used')
24 articles$sovereight[articles$sovereight == 'United States'] <- 'United
    States of America'
25 articles$sovereight[articles$sovereight == 'Tanzania'] <- 'United
    Republic of Tanzania'
26 articles$sovereight[articles$sovereight == 'Serbia'] <- 'Republic of
    Serbia'
27 merged <- merge(World, articles, by='sovereight')
28
29 drop_cnty <- c(
30   "Hong Kong",
31   "Singapore",
32   "Puerto Rico",
33   "Falkland Is.",
34   "Fr. S. Antarctic Lands",
35   "New Caledonia",
36   "Greenland"
```

```

37 )
38 merged <- subset(merged, !(name %in% drop_cnty))
39
40 merged.sp <- as(merged, "Spatial")
41 merged.nb <- poly2nb(pl = merged.sp, queen = TRUE)
42
43 World_centroids <- suppressWarnings(
44   st_centroid(merged) %>%st_coordinates()
45 )
46
47 ##### Country Adjacency #####
48
49 plot(merged.sp, border = "gray")
50 plot(merged.nb, coordinates(merged.sp), col = "red", add = TRUE)
51
52 World.dnb <- dnearneigh(World_centroids, d1 = 0, d2 = 15)
53 plot(merged.sp, border = "gray")
54 plot(World.dnb, World_centroids, col = "brown", add = TRUE)
55
56 World.dnb <- dnearneigh(World_centroids, d1 = 0, d2 = 25)
57 plot(merged.sp, border = "gray")
58 plot(World.dnb, World_centroids, col = "brown", add = TRUE)
59
60 World.knb <- knn2nb(knearneigh(World_centroids, k = 3))
61 plot(merged.sp, border = "gray")
62 plot(World.knb, World_centroids, col = "orange", add = TRUE)
63
64 World.knb <- knn2nb(knearneigh(World_centroids, k = 5))
65 plot(merged.sp, border = "gray")
66 plot(World.knb, World_centroids, col = "orange", add = TRUE)
67
68
69 ##### Local Ratio Mapping #####
70
71 local_ratio_map <- tm_shape(merged)+
72   tm_fill(col = 'local_ratio', title = 'Local Ratio',
73     breaks = c(0, 0.05, 0.1, 0.15, 0.2,0.3, 0.4, 0.6, 1),
74     palette = 'viridis')+
75   tm_layout(main.title = "Local Ratios Calculated per Country",
76     main.title.position = "center")
77
78 local_ratio_map

```

```

79
80
81 ##### Local Ratio Analysis #####
82
83 corr_list <- list()
84 for (m in list(3, 5, 10)) {
85   # Find the m nearest neighbors based on the centroids of each country
86   .
87   World.knb <- knn2nb(knearneigh(World_centroids, k = m))
88
89   # 'country_sma' will hold the estimated local ratios, based on
90   neighbors.
91   country_sma <- list()
92
93   # Iterate over countries.
94   for (i in 1 : nrow(merged)) {
95     neighbor_list <- World.knb[[i]]
96     local_ratios_sum <- 0
97     for (j in 1 : m) {
98       local_ratio <- merged[neighbor_list[[j]],]$local_ratio
99       local_ratios_sum <- local_ratios_sum + local_ratio
100     }
101     country_sma <- append(country_sma, local_ratios_sum/m)
102   }
103
104   # Create data frame for the current m value.
105   country_sma <- as.numeric(country_sma)
106   assign(paste("knb_df_", m, sep = ""), data_frame(merged$sovereight,
107     merged$local_ratio, country_sma))
108 }
109
110 # Rename columns.
111 knb_df_3 <- c('sovereight', 'local_ratio', 'country_sma')
112 knb_df_5 <- c('sovereight', 'local_ratio', 'country_sma')
113 knb_df_10 <- c('sovereight', 'local_ratio', 'country_sma')
114
115 # Compute Pearson correlation coefficients.
116 cor(knb_df_3$local_ratio, knb_df_3$country_sma)
117 cor(knb_df_5$local_ratio, knb_df_5$country_sma)
118 cor(knb_df_10$local_ratio, knb_df_10$country_sma)
119
120 # Merge results into a single data frame that can be shown in the

```

```

    article.
118 knb_df_selected <- data_frame(
119   knb_df_10$sovereight,
120   knb_df_10$local_ratio,
121   knb_df_3$country_sma,
122   knb_df_5$country_sma,
123   knb_df_10$country_sma)
124
125 colnames(knb_df_selected) <- c(
126   'sovereight',
127   'Sample Local Ratio',
128   '3 Neighbor Estimation',
129   '5 Neighbor Estimation',
130   '10 Neighbor Estimation')
131
132 # Create maps for the selected k-neighbors values.
133 empty_world <- World
134 world_with_estimates <- merge(empty_world, knb_df_selected , by = '
    sovereignt', all=T)
135
136 three_neighbor_map <- tm_shape(world_with_estimates)+
137   tm_fill(col = '3 Neighbor Estimation', title = '3 Neighbor Ratio
    Estimation',
138     breaks = c(0, 0.05, 0.1, 0.15, 0.2, 0.3, 0.4, 0.6, 1),
139     palette = 'viridis')
140 three_neighbor_map
141
142 five_neighbor_map <- tm_shape(world_with_estimates)+
143   tm_fill(col = '5 Neighbor Estimation', title = '5 Neighbor Ratio
    Estimation',
144     breaks = c(0, 0.05, 0.1, 0.15, 0.2, 0.3, 0.4, 0.6, 1),
145     palette = 'viridis')
146 five_neighbor_map
147
148 ten_neighbor_map <- tm_shape(world_with_estimates)+
149   tm_fill(col = '10 Neighbor Estimation', title = '10 Neighbor Ratio
    Estimation',
150     breaks = c(0, 0.05, 0.1, 0.15, 0.2, 0.3, 0.4, 0.6, 1),
151     palette = 'viridis')
152 ten_neighbor_map
153
154 # Compute Moran I test statistics

```

```
155
156 # 3 neighbors: simple average
157 merged.knb_for_moran_3 <- knn2nb(knearneigh(World_centroids, k = 3))
158 merged.w3 <- nb2listw(merged.knb_for_moran_3)
159 moran.test(merged$local_ratio, merged.w3)
160
161 # 3 neighbors: IDW
162 merged.idw_w3 <- nb2listwdist(merged.knb_for_moran_3, merged)
163 moran.test(merged$local_ratio, merged.idw_w3)
164
165 # 5 neighbors: simple average
166 merged.knb_for_moran_5 <- knn2nb(knearneigh(World_centroids, k = 5))
167 merged.w5 <- nb2listw(merged.knb_for_moran_5)
168 moran.test(merged$local_ratio, merged.w5)
169
170 # 5 neighbors: IDW
171 merged.idw_w5 <- nb2listwdist(merged.knb_for_moran_5, merged)
172 moran.test(merged$local_ratio, merged.idw_w5)
173
174 # 10 neighbors: simple average
175 merged.knb_for_moran_10 <- knn2nb(knearneigh(World_centroids, k = 10))
176 merged.w10 <- nb2listw(merged.knb_for_moran_10)
177 moran.test(merged$local_ratio, merged.w10)
178
179 # 10 neighbors: IDW
180 merged.idw_w10 <- nb2listwdist(merged.knb_for_moran_10, merged)
181 moran.test(merged$local_ratio, merged.idw_w10)
```

Listing 3: R Visualization and Analysis of Local Ratios