

Final Project

CS110: Computation: Solving Problems with Algorithms

Kalia Barkai

December 2017

R-Trees: Introduction

R-Trees are a type of data structure used to store and search spatial indices. The advantage of searching through R-Trees for values, is that it is a multi-dimensional data structure and therefore, can find the values within a certain spatial area - and the values near it.

It does this by storing the key, along with it's tuple coordinates, in rectangles (i.e. areas), which are child nodes of greater rectangles (i.e. the encompassing areas). The tree structure is similar to that of a B+-Tree, where it is balanced by having equal lengths of the paths from the root to all leaf nodes, and each node can have 0 to N children nodes. N is the decided upon max number of sibling nodes and this applies to the root level too.

The following are the invariants of an R-Tree:

1. Every node (non-leaf) has between n and N children (except for root), where $n \leq \frac{N}{2}$
2. The root has at least two children, unless it is a leaf
3. Each node rectangle is the smallest rectangle that contains its children
4. All leaves are on the same level (balanced tree)

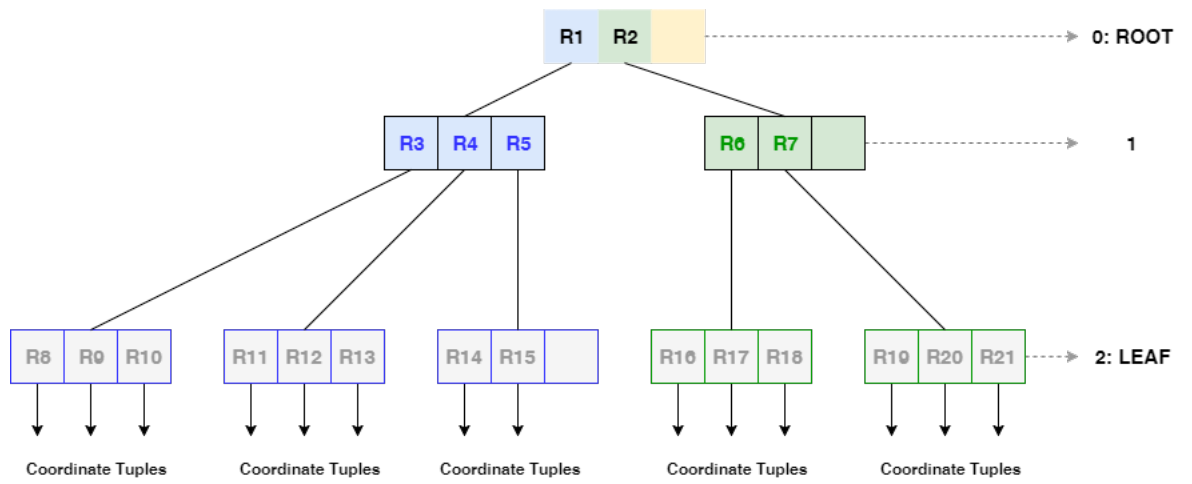


Figure 1. R-Tree example with 21 rectangles. $N = 3$

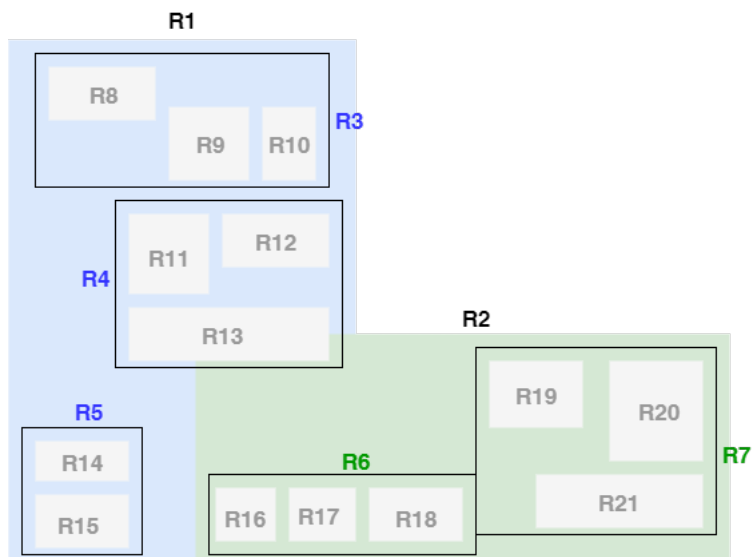


Figure 2. R-Tree rectangles example from R-Tree in Figure 1.

The R-Tree datastructure supports the following operations: insert, search and deletion. Search can be implemented in two more ways: range search or priority search. Range search finds all the values within a rectangle/area, while priority search finds a specific value. Priority search might be used for a nearest neighbour search for example.

Range search happens recursively where any rectangle that intersects with the search rectangle is searched until it reaches the leaf nodes, which are then checked against the main search rectangle.

Insertion also happens recursively from the root node. Each rectangle is checked by seeing which rectangle would need the least enlargement for the insertion to take place, it then goes into this rectangle checking the child nodes in the same way until it reaches a leaf node where, if there is space, the value is inserted in the leaf node, and if there is no space, then the leaf node is split.

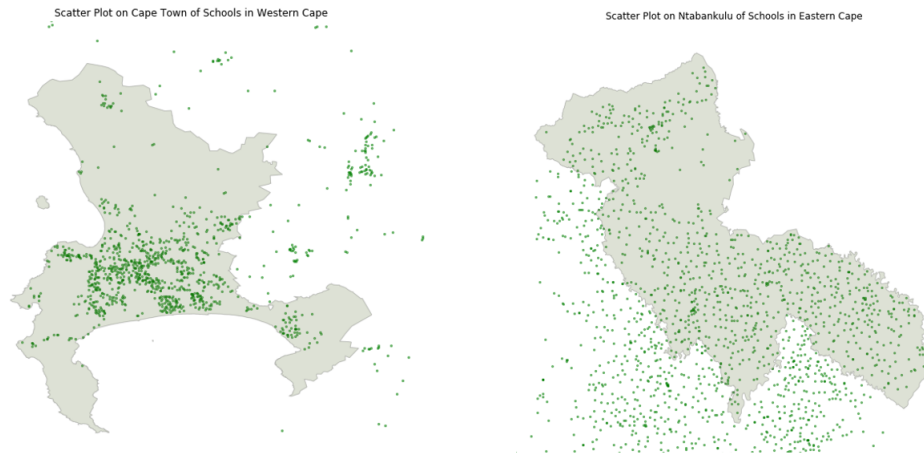
Deletion either updates the bounding boxes of parent rectangles in order to rebalance, however if the R-Tree is not full, then it will be deleted and each value will be reinserted into the tree.

The average case time complexity for search in an R-Tree is $T(n) = O(\log n)$, where n refers to the number of rectangles/nodes. This assumes that there is not too much overlap in the tree and only a minimal amount of paths must be followed for the search.

R-Trees: Application

Using the schools database from the Department of Basic Education of the Republic of South Africa which is divided up by provinces, Cape Town and Ntabankulu student to teacher ratios at each school were compared.¹

Since the full dataset is divided by provinces and not by cities, but includes coordinates for each school, it is possible to use spatial indexing in R-Trees to find which schools fall within chosen cities.



Figures 3. and 4. Scatter plot of schools in Cape Town and in Ntabankulu and surrounding area. Created using Geopandas and Shapely libraries in Python. Schools database from the Department of Basic Education, RSA (2016).

Listing 1. below shows the code which implements search in the R-Tree of the Western Cape school nodes.²

```
1 # network of schools coordinates
2 gdf_nodes = gpd.GeoDataFrame(data={'x':x, 'y':y})
3 # system of starting coordinates (dict)
4 gdf_nodes.crs = ct.crs
```

¹Cape Town is in the top 3 richest cities in South Africa while Ntabankulu (part of the Alfred Nzo Municipality) is in the poorest 3 cities.

²The rest of the code can be found in the Appendix.

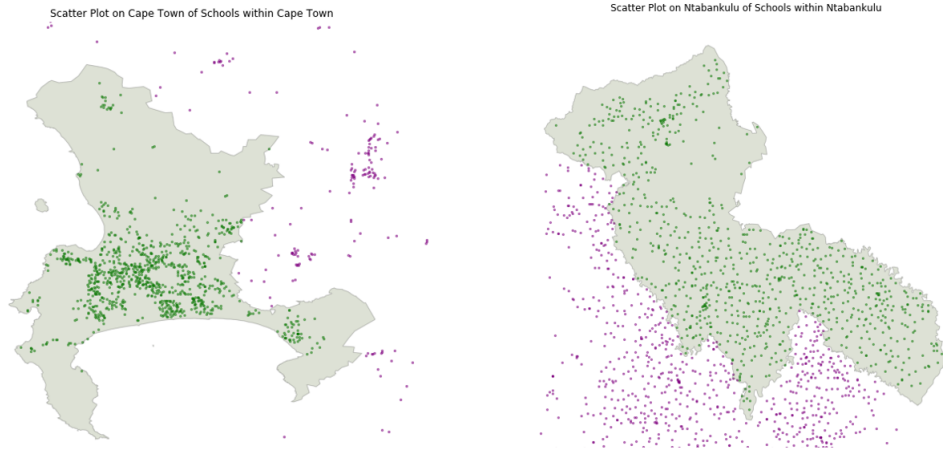
```

5 # value is the ratio of learners per educator
6 gdf_nodes.name = lpe
7 # create 'geometry' col with coordinates
8 gdf_nodes['geometry'] = gdf_nodes.apply(lambda row:
    Point((row['x'], row['y'])), axis=1)
9
10 spatial_index = gdf_nodes.sindex # spatial index
11 # check which spatial indices intersect with our
    bounding box of the polygon (our city)
12 # possible, because there are no false negatives,
    but there can be false positives
13 # this happens when the bounding box in the R-Tree
    of the polygon has bounds greater than the
    polygon bounds
14 # in this case points outside the city but in this
    bounding box are false positives
15 possible_index = list(spatial_index.intersection(
    geometry.bounds))
16 # fetch nodes of possible schools
17 possible_schools = gdf_nodes.iloc[possible_index]
18 # check if the possible schools intersect with our
    actual polygon
19 # this gives us the actual schools within
20 schools_ct = possible_schools[possible_schools.
    intersects(geometry)]

```

Listing 1: Code which implements spatial indexing with R-Trees using the Geopandas Python library.

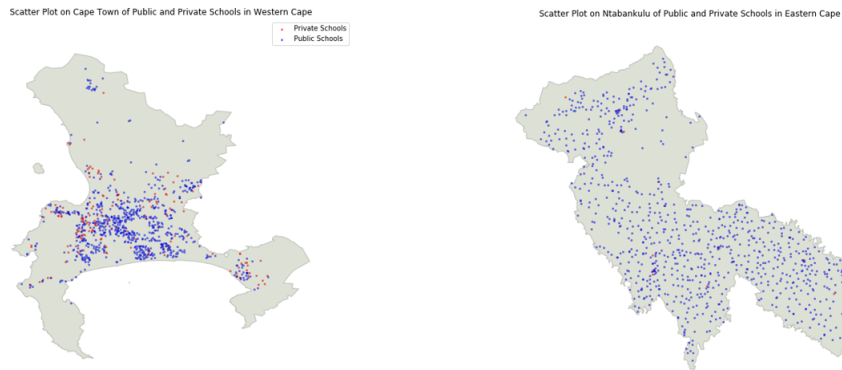
Figures 5. and 6. show the visualisations created by finding which schools lie within the cities and which do not.



Figures 5. and 6. Scatter plot of schools within Cape Town and in Ntabankulu and surrounding area. R-Tree spatial indexing used to find points within the cities. Created using Geopandas and Shapely libraries in Python. Schools database from the Department of Basic Education, RSA (2016).

Using the values of the schools within the cities, we can calculate whether there is a significant difference between the amount of students per educators at each school. First, the mean for schools in Cape Town was calculated to be 30.055, while the mean for schools in Ntabankulu was calculated to be 32.622. However, schools in Cape Town have a range of 3 to 67 students per educator, while schools in Ntabankulu have a range of 5 to 83 students per educator. To calculate whether there is actually a significant difference, it is important to look at the p-value between the two cities. Using scipy.stats library in Python, the p-value = 0.00000043. Since this is ≤ 0.05 , there is a significant difference. Therefore, only looking at the mean of this data, might be misleading.

For further analysis we can also look at a visualisation of the private and public schools within each city. This can be seen in Figures 7. and 8. below.



Figures 7. and 8. Scatter plot of public and private schools in Cape Town and in Ntabankulu and surrounding area. Created using Geopandas and Shapely libraries in Python. Schools database from the Department of Basic Education, RSA (2016).

As seen in Figures 7. and 8., Cape Town has more private schools. This is just one further analysis of the two regions. More analysis on the populations of the areas, the schools budget and resources and the socio-economic disposition of the communities surrounding the schools would be relevant to further understand the reason for the significant difference between students per educator in each region.

Appendix

Notebook of all code available on [here](#).

```
1 import osmnx as ox, matplotlib.pyplot as plt, pandas
    as pd, geopandas as gpd
2 from descartes import PolygonPatch
3 from shapely.geometry import Point, Polygon,
    MultiPolygon
4 %matplotlib inline
5 # for high quality visualisations:
6 %config InlineBackend.figure_format = 'svg'
7
8 # boundary of Cape Town
9 ct = ox.gdf_from_place('City of Cape Town, South
    Africa')
10
11 # x and y bounds
12 west, south, east, north = ct.unary_union.buffer
    (0.1).bounds
13
14 # geometry into polygon
15 geometry = ct['geometry'].iloc[0]
16 if isinstance(geometry, Polygon):
17     geometry = MultiPolygon([geometry])
18
19 # Western Cape school database
20 wc_schools = pd.read_excel("wc.xlsx")
```



```

21
22 # Latitude coordinate of school
23 y = wc_schools["GIS_Latitude"]
24 # Longitude coordinate of school
25 x = wc_schools["GIS_Longitude"]
26
27 # amount of learners in school
28 learners = wc_schools["Learners_2016"]
29 # amount of educators in school
30 educators = wc_schools["Educators_2016"]
31 lpe = [] # learners per educator
32 for i in range(len(learners)):
33     lpe.append(learners[i]/educators[i])
34
35 gdf_nodes = gpd.GeoDataFrame(data={'x':x, 'y':y})
36 gdf_nodes.crs = ct.crs
37 gdf_nodes.name = lpe
38 gdf_nodes['geometry'] = gdf_nodes.apply(lambda row:
39     Point((row['x'], row['y'])), axis=1)
40
41 fig, ax = plt.subplots(figsize=(10,10))
42 for polygon in geometry:
43     patch = PolygonPatch(polygon, fc='#556B2F', ec='
44     k', alpha=0.2, zorder=2)
45     ax.add_patch(patch)

```

```

45 plt.title("Scatter Plot on Cape Town of Schools in
    Western Cape")
46 ax.scatter(x=gdf_nodes['x'], y=gdf_nodes['y'], s=5,
    c='g', zorder=1, alpha = 0.5) # plotting schools
47 ax.set_xlim(west, east)
48 ax.set_ylim(south, north)
49 ax.axis('off')
50 plt.savefig("ct_1.png")
51 plt.show()

```

Listing 2: Code used to create Figure 3

```

1 # Compare to Ntabankulu in Eastern Cape
2 # repeat process
3 nt = ox.gdf_from_place('Alfred Nzo District
    Municipality, South Africa')
4
5 # Eastern Cape schools database
6 ec = pd.read_excel("ec.xlsx")
7
8 y_ec = ec["GIS_Latitude"]
9 x_ec = ec["GIS_Longitude"]
10
11 learners_ec = ec["Learners_2016"]
12 educators_ec = ec["Educators_2016"]
13 lpe_ec = []
14 for i in range(len(learners_ec)):
15     lpe_ec.append(learners_ec[i]/educators_ec[i])

```

```

16
17 geo_nt = nt['geometry'].iloc[0]
18 if isinstance(geo_nt, Polygon):
19     geo_nt = MultiPolygon([geo_nt])
20
21 gdf_nodes_ec = gpd.GeoDataFrame(data={'x':x_ec, 'y':
    y_ec})
22 gdf_nodes_ec.crs = nt.crs
23 # value is the ratio of learners per educator
24 gdf_nodes_ec.name = lpe_ec
25 gdf_nodes_ec['geometry'] = gdf_nodes_ec.apply(lambda
    row: Point((row['x'], row['y'])), axis=1)
26
27 west_nt, south_nt, east_nt, north_nt = nt.
    unary_union.buffer(0.1).bounds
28
29 fig, ax = plt.subplots(figsize=(10,10))
30 for polygon in geo_nt:
31     patch = PolygonPatch(polygon, fc='#556B2F', ec='
    k', alpha=0.2, zorder=2)
32     ax.add_patch(patch)
33
34
35 plt.title("Scatter Plot on Ntabankulu of Schools in
    Eastern Cape")
36 ax.scatter(x=gdf_nodes_ec['x'], y=gdf_nodes_ec['y'],
    s=5, c='g', zorder=1, alpha = 0.5)

```

```

37 ax.set_xlim(west_nt, east_nt)
38 ax.set_ylim(south_nt, north_nt)
39 ax.axis('off')
40 plt.savefig("nt_1.png")
41 plt.show()

```

Listing 3: Code used to create Figure 4

```

1 # schools outside Cape Town
2 schools_not = gdf_nodes[~gdf_nodes.isin(schools_ct)]
3
4 fig, ax = plt.subplots(figsize=(10,10))
5 for polygon in geometry:
6     patch = PolygonPatch(polygon, fc='#556B2F', ec='
7     k', alpha=0.2, zorder=2)
8     ax.add_patch(patch)
9
10 plt.title("Scatter Plot on Cape Town of Schools
11     within Cape Town")
12
13 ax.scatter(x=schools_not['x'], y=schools_not['y'], s
14     =5, c='purple', zorder=1, alpha = 0.5)
15 ax.scatter(x=schools_ct['x'], y=schools_ct['y'], s
16     =5, c='g', zorder=2, alpha = 0.5)
17
18 ax.set_xlim(west, east)
19 ax.set_ylim(south, north)
20 ax.axis('off')
21 plt.savefig("ct_3.png")

```

```
17 plt.show()
```

Listing 4: Code used to create Figure 5

```
1 spatial_index_ec = gdf_nodes_ec.sindex
2 possible_index_ec = list(spatial_index_ec.
    intersection(geo_nt.bounds))
3 possible_schools_ec = gdf_nodes_ec.iloc[
    possible_index_ec]
4 schools_nt = possible_schools_ec[possible_schools_ec
    .intersects(geo_nt)]
5
6 schools_not_nt = gdf_nodes_ec[~gdf_nodes_ec.isin(
    schools_nt)]
7
8 fig, ax = plt.subplots(figsize=(10,10))
9 for polygon in geo_nt:
10     patch = PolygonPatch(polygon, fc='#556B2F', ec='
    k', alpha=0.2, zorder=2)
11     ax.add_patch(patch)
12
13 ax.scatter(x=schools_not_nt['x'], y=schools_not_nt['
    y'], s=5, c='purple', zorder=1, alpha = 0.5)
14 ax.scatter(x=schools_nt['x'], y=schools_nt['y'], s
    =5, c='g', zorder=2, alpha = 0.5)
15
16 plt.title("Scatter Plot on Ntabankulu of Schools
    within Ntabankulu")
```

```

17 ax.set_xlim(west_nt, east_nt)
18 ax.set_ylim(south_nt, north_nt)
19 ax.axis('off')
20 plt.savefig("nt_3.png")
21 plt.show()

```

Listing 5: Code used to create Figure 6

```

1 # public or independent school
2 sector = wc_schools["Sector"]
3
4 private_lon = []
5 private_lat = []
6 public_lon = []
7 public_lat = []
8
9 for i in schools_ct.index:
10     # private school coordinates
11     if sector[i] == "INDEPENDENT":
12         private_lon.append(x[i])
13         private_lat.append(y[i])
14     else: # public school coordinates
15         public_lon.append(x[i])
16         public_lat.append(y[i])
17
18
19 fig, ax = plt.subplots(figsize=(10,10))
20 for polygon in geometry:

```

```

21     patch = PolygonPatch(polygon, fc='#556B2F', ec='
      k', alpha=0.2, zorder=2)
22     ax.add_patch(patch)
23
24 plt.title("Scatter Plot on Cape Town of Public and
      Private Schools in Western Cape")
25 ax.scatter(private_lon, private_lat, s=5, c='r',
      zorder=1, alpha = 0.5, label = "Private Schools")
26 ax.scatter(public_lon, public_lat, s=5, c='b',
      zorder=1, alpha = 0.5, label = "Public Schools")
27 ax.legend()
28
29 ax.set_xlim(west, east)
30 ax.set_ylim(south, north)
31 ax.axis('off')
32 plt.savefig("ct_2.png")
33 plt.show()

```

Listing 6: Code used to create Figure 7

```

1 # public or independent school
2 sector_ec = ec["Sector"]
3
4 private_lon_ec = []
5 private_lat_ec = []
6 public_lon_ec = []
7 public_lat_ec = []
8

```

```

9 for i in schools_nt.index:
10     if sector_ec[i] == "INDEPENDENT":
11         private_lon_ec.append(x_ec[i])
12         private_lat_ec.append(y_ec[i])
13     else:
14         public_lon_ec.append(x_ec[i])
15         public_lat_ec.append(y_ec[i])
16
17
18 fig, ax = plt.subplots(figsize=(10,10))
19 for polygon in geo_nt:
20     patch = PolygonPatch(polygon, fc='#556B2F', ec='
k', alpha=0.2, zorder=2)
21     ax.add_patch(patch)
22
23 plt.title("Scatter Plot on Ntabankulu of Public and
Private Schools in Eastern Cape")
24 ax.scatter(private_lon_ec, private_lat_ec, s=5, c='r
', zorder=1, alpha = 0.5, label = "Private
Schools")
25 ax.scatter(public_lon_ec, public_lat_ec, s=5, c='b',
zorder=1, alpha = 0.5, label = "Public Schools")
26 ax.legend()
27
28 ax.set_xlim(west_nt, east_nt)
29 ax.set_ylim(south_nt, north_nt)
30 ax.axis('off')

```



```
31 plt.savefig("nt_2.png")
32 plt.show()
```

Listing 7: Code used to create Figure 8

```
1 # empty list of teacher-student ratio for schools in
   Cape Town
2 schools_tsr = []
3 for i in range(len(schools_ct)):
4     index = schools_ct.index[i]
5     schools_tsr.append(gdf_nodes.name[index])
6
7 schools_tsr_nt = []
8 for i in range(len(schools_nt)):
9     index = schools_nt.index[i]
10    schools_tsr_nt.append(gdf_nodes_ec.name[index])
11
12 from scipy.stats import ttest_ind
13 t, p = ttest_ind(tsr, tsr_nt)
14 print p
15 >> 4.27283063584e-07
```

Listing 8: Code used to calculate p-value of ratio of students per educator in Cape Town versus Ntabankulu.

References

- “A Short Tutorial On B+ Tree”. 2015. Sketching Dream. <http://sketchingdream.com/blog/b-plus-tree-tutorial/>. “B+ Tree”. 2017. Wikipedia. https://en.wikipedia.org/wiki/B%2B_tree.
- Boeing, Geoff. 2016. “R-Tree Spatial Indexing With Python - Geoff Boeing”. Geoff Boeing. <http://geoffboeing.com/2016/10/r-tree-spatial-index-python/>.
- Boeing, Geoff. 2017. “Rtree-Spatial-Indexing”. Github. <https://github.com/gboeing/urban-data-science/blob/master/19-Spatial-Analysis-and-Cartography/r-tree-spatial-indexing.ipynb>.
- “EMIS Downloads”. 2017. Education.Gov.Za. <https://www.education.gov.za/Programmes/EMIS/EMISDownloads.aspx>. Guttman, Antonin. 1984. “R-Trees”. ACM SIGMOD Record 14 (2): 47. doi:10.1145/971697.602266.
- “R-Tree”. 2017. Wikipedia. <https://en.wikipedia.org/wiki/R-tree>. R-Tree Chapter 1. 2008. Ebook. Itoma. <http://www.bowdoin.edu/ltoma/teaching/cs340/spring08/Papers/Rtree-chap1.pdf>.