

# Data exploration with GeoPandas

**CS424: Visualization & Visual Analytics**

**Fabio Miranda**

**<https://fmiranda.me>**

# GeoPandas



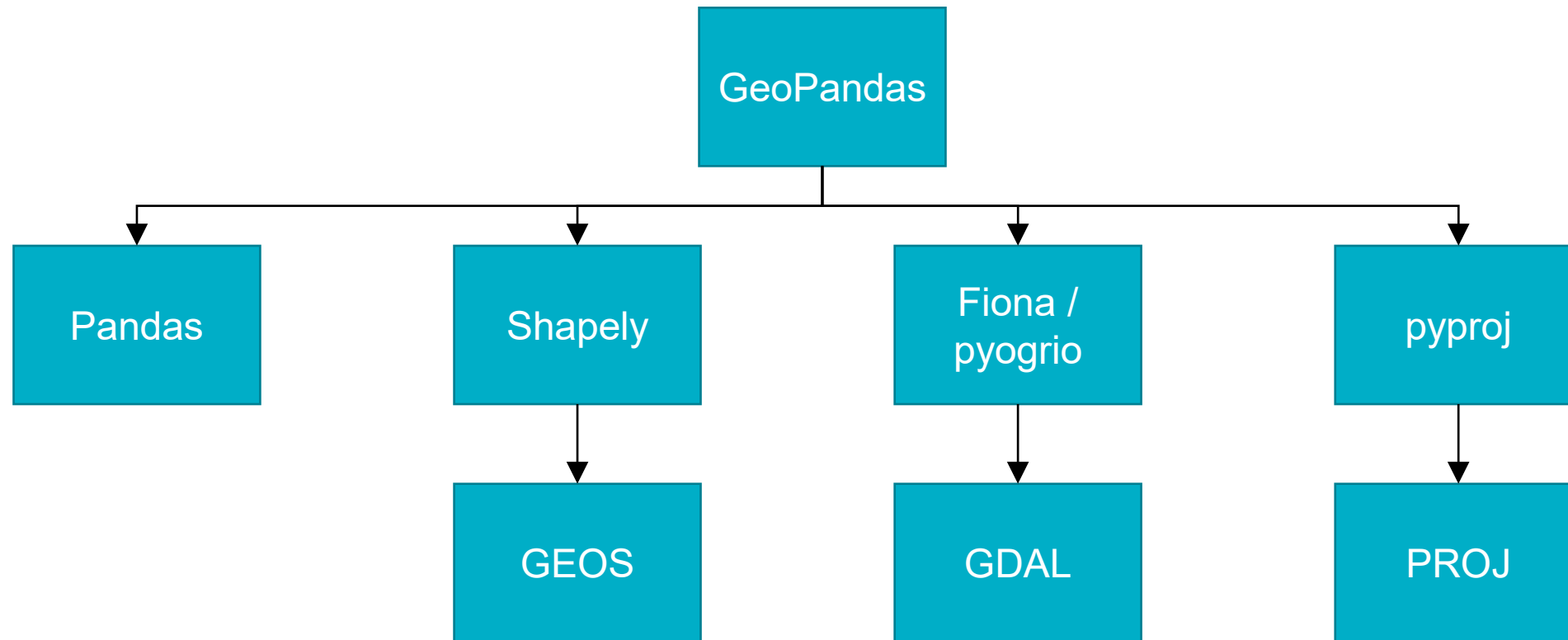
- Geospatial data & Python made easier.
- Extends Pandas to work with geographic objects and spatial operations.
- Combines the power of several libraries (Pandas, geos, shapely, gdal, pyproj, rtree, ...)

# GeoPandas



- Read and write several geo formats (Fiona, GDAL).
- Usual DataFrame manipulation.
- Element-wise spatial operations (intersection, union, difference, ...)
- Re-project data.
- Visualize geometries.
- Advanced spatial operations: spatial joins and overlays.

# GeoPandas



# Shapely

- Python package for the manipulation of geometric objects.

```
1 from shapely.geometry import Point, LineString, Polygon
2
3 point = Point(1, 1)
4 line = LineString([(0, 0), (1, 2), (2, 2)])
5 poly = line.buffer(1)
```

```
1 line
```



```
1 poly
```

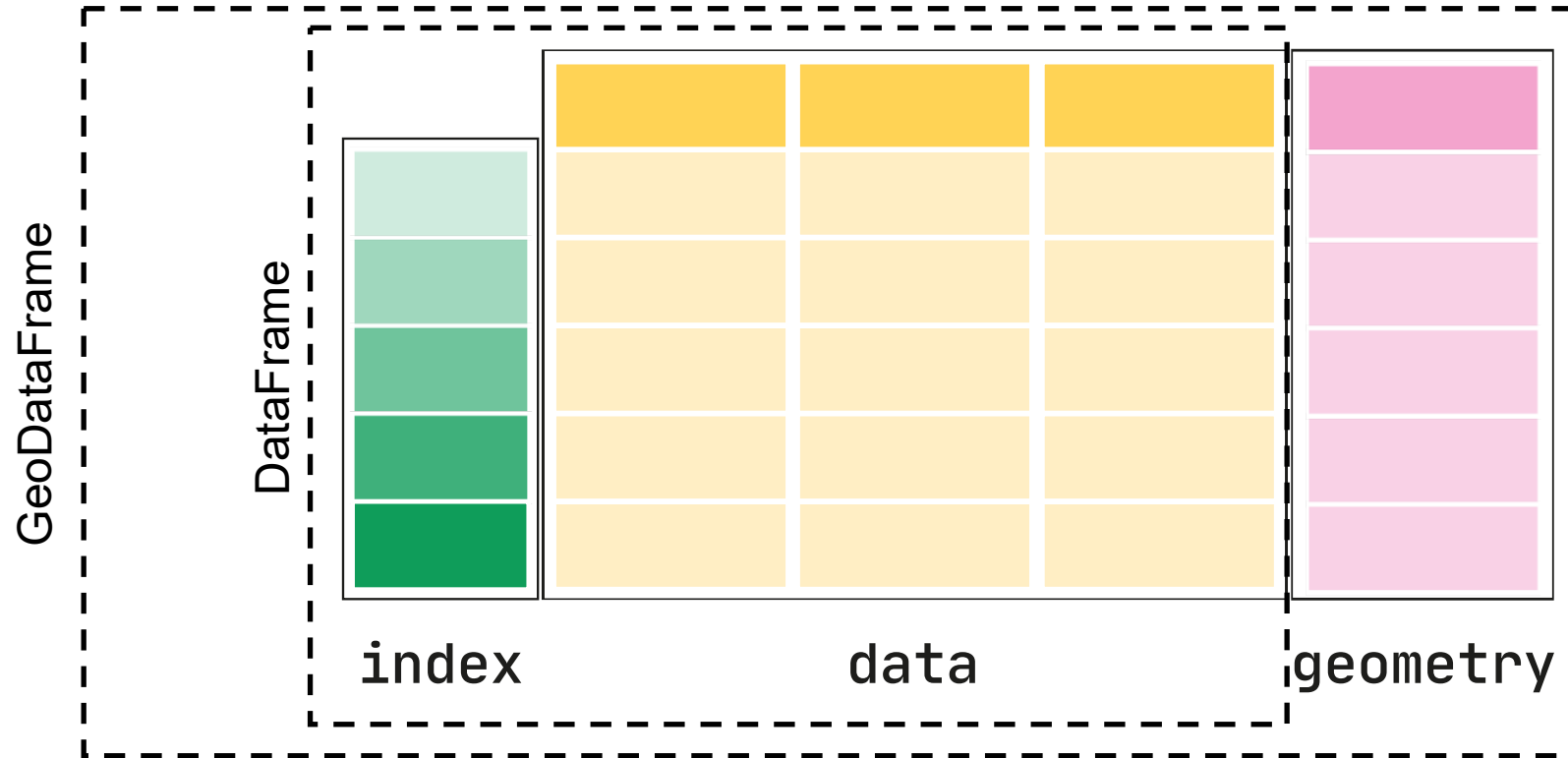


```
1 poly.contains(point)
```

True

# GeoPandas

- Core data structure in GeoPandas is the GeoDataFrame (subclass of Pandas' DataFrame).
- It can store geometry columns and perform spatial operations.

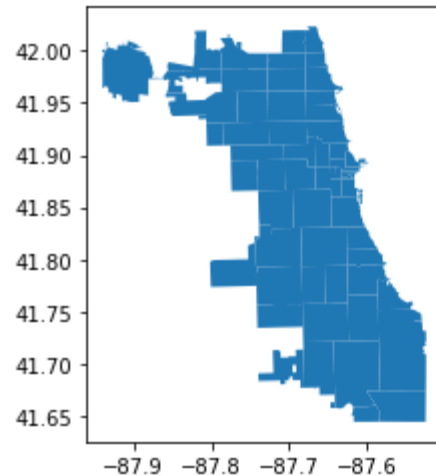


# Reading and writing files

```
1 gdf = gpd.read_file('chicago.geojson')
```

```
1 gdf.plot()
```

<AxesSubplot:>



```
1 gdf
```

	objectid	shape_area	shape_len	zip	geometry
0	33	106052287.488	42720.0444058	60647	MULTIPOLYGON (((-87.67762 41.91776, -87.67761 ...
1	34	127476050.762	48103.7827213	60639	MULTIPOLYGON (((-87.72683 41.92265, -87.72693 ...
2	35	45069038.4783	27288.6096123	60707	MULTIPOLYGON (((-87.78500 41.90915, -87.78531 ...
3	36	70853834.3797	42527.9896789	60622	MULTIPOLYGON (((-87.66707 41.88885, -87.66707 ...
4	37	99039621.2518	47970.1401531	60651	MULTIPOLYGON (((-87.70656 41.89555, -87.70672 ...
...	...	...	...	...	...
56	57	155285532.005	53406.9156168	60623	MULTIPOLYGON (((-87.69479 41.83008, -87.69486 ...
57	58	211114779.439	58701.3253749	60629	MULTIPOLYGON (((-87.68306 41.75786, -87.68306 ...
58	59	211696050.967	58466.1602979	60620	MULTIPOLYGON (((-87.62373 41.72167, -87.62388 ...
59	60	125424284.172	52377.8545408	60637	MULTIPOLYGON (((-87.57691 41.79511, -87.57700 ...
60	61	167872012.644	53040.9070778	60619	MULTIPOLYGON (((-87.58592 41.75150, -87.58592 ...

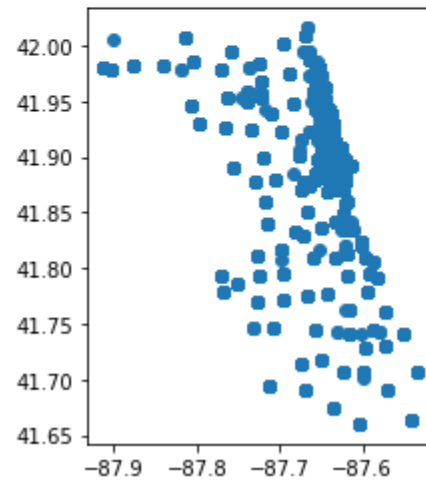
61 rows × 5 columns

# Reading CSV file

```
1 df = pd.read_csv('data/Taxi_Trips.csv')
2 geometry = [Point(xy) for xy in zip(df['Pickup Centroid Longitude'], df['Pickup Centroid Latitude'])]
3 gdf = gpd.GeoDataFrame(df, geometry=geometry, crs=4326)
```

```
1 gdf.plot()
```

<AxesSubplot:>



Geometry from lon / lat (x / y)



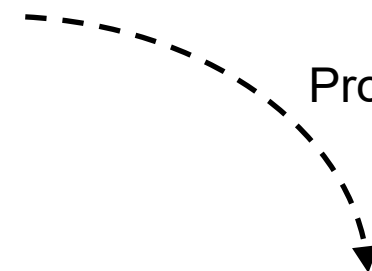
# Projections

```
1 gdf
```

	objectid	shape_area	shape_len	zip	geometry
0	33	106052287.488	42720.0444058	60647	MULTIPOLYGON (((-87.67762 41.91776, -87.67761 ...
1	34	127476050.762	48103.7827213	60639	MULTIPOLYGON (((-87.72683 41.92265, -87.72693 ...
2	35	45069038.4783	27288.6096123	60707	MULTIPOLYGON (((-87.78500 41.90915, -87.78531 ...
3	36	70853834.3797	42527.9896789	60622	MULTIPOLYGON (((-87.66707 41.88885, -87.66707 ...
4	37	99039621.2518	47970.1401531	60651	MULTIPOLYGON (((-87.70656 41.89555, -87.70672 ...
...	...	...	...	...	...
56	57	155285532.005	53406.9156168	60623	MULTIPOLYGON (((-87.69479 41.83008, -87.69486 ...
57	58	211114779.439	58701.3253749	60629	MULTIPOLYGON (((-87.68306 41.75786, -87.68306 ...
58	59	211696050.967	58466.1602979	60620	MULTIPOLYGON (((-87.62373 41.72167, -87.62388 ...
59	60	125424284.172	52377.8545408	60637	MULTIPOLYGON (((-87.57691 41.79511, -87.57700 ...
60	61	167872012.644	53040.9070778	60619	MULTIPOLYGON (((-87.58592 41.75150, -87.58592 ...

61 rows × 5 columns

Projecting to EPSG:3395



```
1 gdf = gdf.to_crs("EPSG:3395")
```

```
1 gdf
```

	objectid	shape_area	shape_len	zip	geometry
0	33	106052287.488	42720.0444058	60647	MULTIPOLYGON (((-9760228.181 5120114.708, -976...
1	34	127476050.762	48103.7827213	60639	MULTIPOLYGON (((-9765706.326 5120843.341, -976...
2	35	45069038.4783	27288.6096123	60707	MULTIPOLYGON (((-9772181.764 5118831.519, -977...
3	36	70853834.3797	42527.9896789	60622	MULTIPOLYGON (((-9759053.446 5115807.386, -975...
4	37	99039621.2518	47970.1401531	60651	MULTIPOLYGON (((-9763449.188 5116805.817, -976...
...	...	...	...	...	...
56	57	155285532.005	53406.9156168	60623	MULTIPOLYGON (((-9762139.685 5107055.000, -976...
57	58	211114779.439	58701.3253749	60629	MULTIPOLYGON (((-9760833.554 5096312.536, -976...
58	59	211696050.967	58466.1602979	60620	MULTIPOLYGON (((-9754228.915 5090933.835, -975...
59	60	125424284.172	52377.8545408	60637	MULTIPOLYGON (((-9749017.532 5101851.550, -974...
60	61	167872012.644	53040.9070778	60619	MULTIPOLYGON (((-9750019.820 5095367.709, -975...

61 rows × 5 columns

# Projections

```
1 | gdf
```

	objectid	shape_area	shape_len	zip	geometry
0	33	108			<Geographic 2D CRS: EPSG:4326>
1	34	127			Name: WGS 84
2	35	450			Axis Info [ellipsoidal]:
3	36	708			- Lat[north]: Geodetic latitude (degree)
4	37	990			- Lon[east]: Geodetic longitude (degree)
...	...	...			Area of Use:
56	57	155			- name: World.
57	58	211			- bounds: (-180.0, -90.0, 180.0, 90.0)
58	59	211			Datum: World Geodetic System 1984 ensemble
59	60	125424284.172	52377.8545408	60637	- Ellipsoid: WGS 84
60	61	167872012.644	53040.9070778	60619	- Prime Meridian: Greenwich

61 rows × 5 columns

Projecting to EPSG:3395

```
1 | gdf = gdf.to_crs("EPSG:3395")
```

```
1 | gdf
```

	objectid	shape_area	shape_len	zip	geometry
0	33	1060522			<Derived Projected CRS: EPSG:3395>
1	34	1274760			Name: WGS 84 / World Mercator
2	35	4506903			Axis Info [cartesian]:
3	36	7085383			- E[east]: Easting (metre)
4	37	9903962			- N[north]: Northing (metre)
...	...	...			Area of Use:
56	57	1552855			- name: World between 80°S and 84°N.
57	58	2111147			- bounds: (-180.0, -80.0, 180.0, 84.0)
58	59	2116960			Coordinate Operation:
59	60	1254242			- name: World Mercator
60	61	167872012.644	53040.9070778	60619	- method: Mercator (variant A)

61 rows × 5 columns

# Simple operations

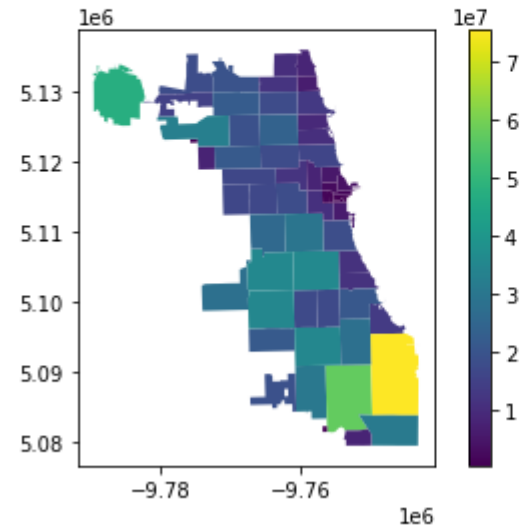
- Accessing & plotting geometry area

```
1 gdf.area
0    1.774279e+07
1    2.132685e+07
2    7.540024e+06
3    1.184732e+07
4    1.656003e+07
...
56   2.592093e+07
57   3.515998e+07
58   3.521726e+07
59   2.089192e+07
60   2.793029e+07
Length: 61, dtype: float64
```

```
1 gdf['area'] = gdf.area
```

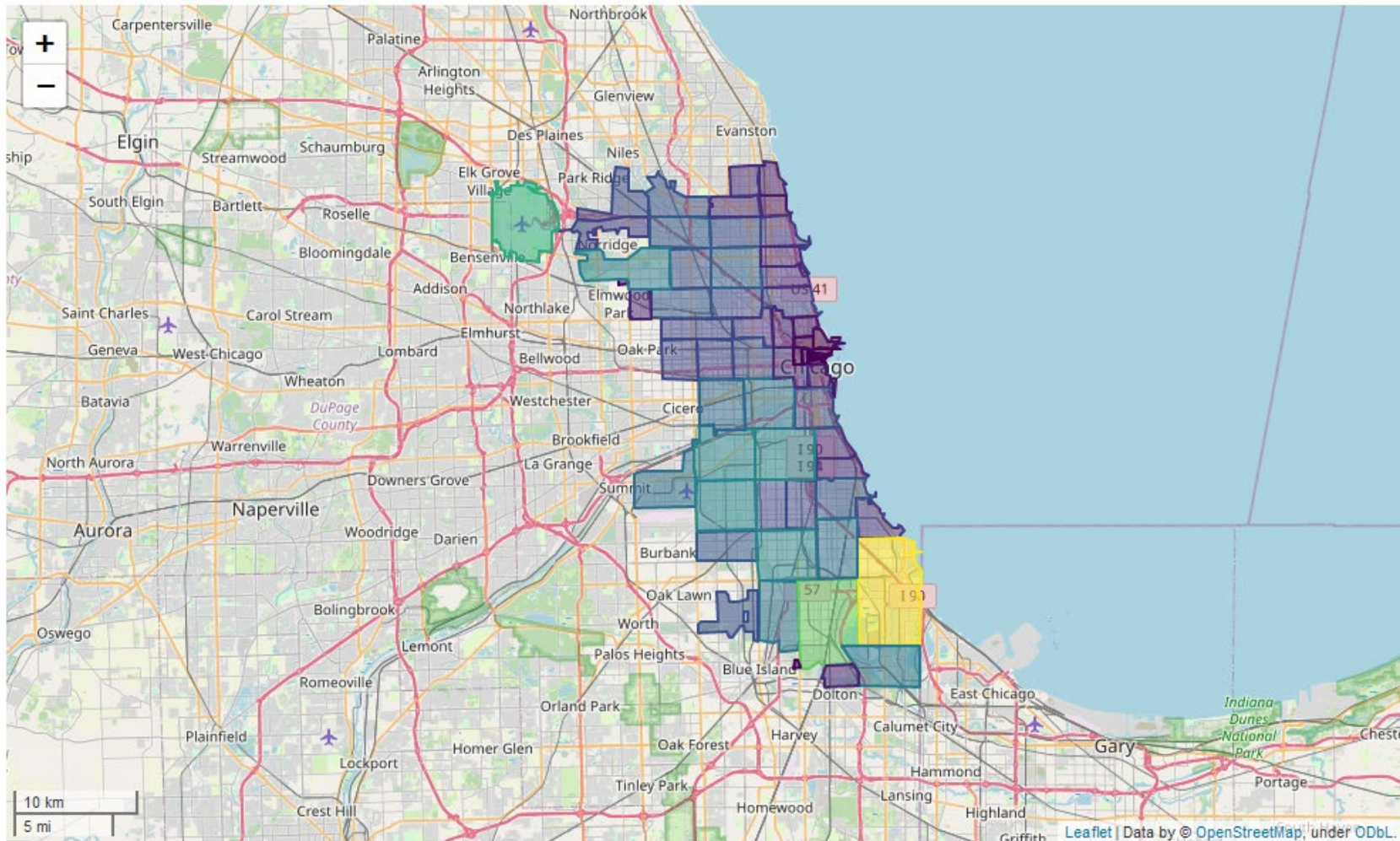
```
1 gdf.plot("area", legend=True)
```

<AxesSubplot:>



# Simple operations

```
1 gdf.explore("area", legend=False)
```



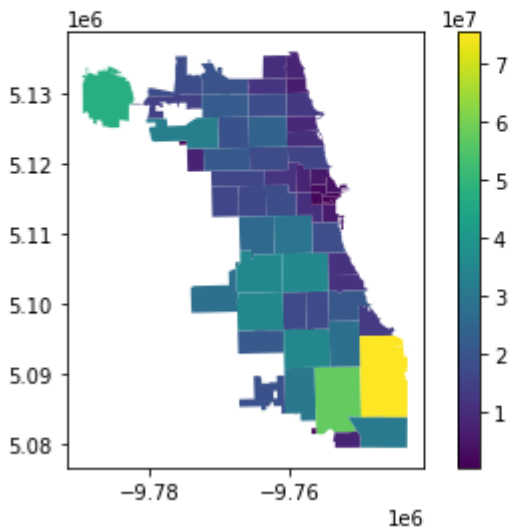
# Simple operations

Changing geometry

```
1 gdf['area'] = gdf.area
```

```
1 gdf.plot("area", legend=True)
```

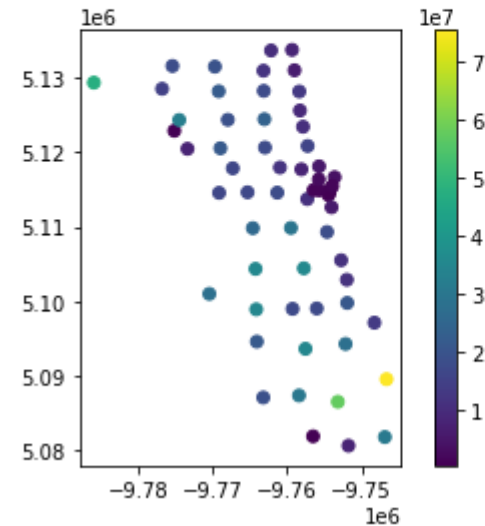
<AxesSubplot:>



```
1 gdf['centroid'] = gdf.centroid
```

```
1 gdf = gdf.set_geometry('centroid')  
2 gdf.plot("area", legend=True)
```

<AxesSubplot:>



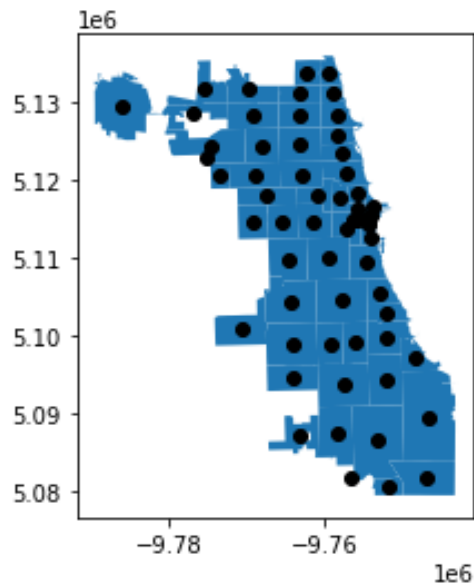
## Simple operations

- Plotting centroids and zip areas:

```
1 gdf = gdf.set_geometry('geometry')  
2 ax = gdf['geometry'].plot()  
3 gdf['centroid'].plot(ax=ax, color="black")
```

Setting geometry back

<AxesSubplot:>





# Geometry operations

```
1 gdf["convex_hull"] = gdf.convex_hull
```

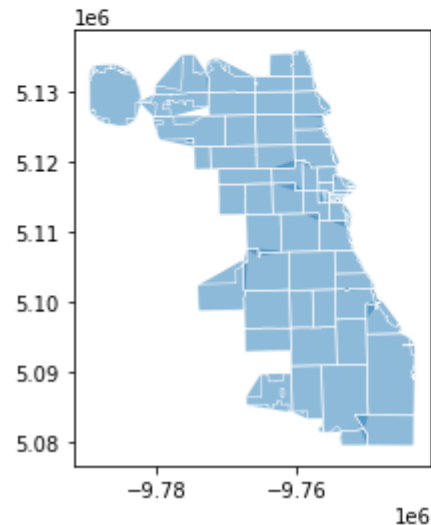
→ New column with convex hull

```
1 gdf['boundary'] = gdf.boundary
```

→ New column with boundary (i.e., outline)

```
1 ax = gdf["convex_hull"].plot(alpha=.5)
2 gdf['boundary'].plot(ax=ax, color="white", linewidth=.5)
```

<AxesSubplot:>



# Geometry operations

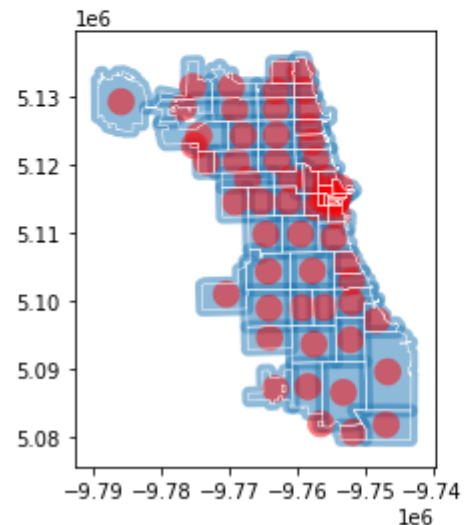
```
1 gdf["buffered"] = gdf.buffer(1000)
2 gdf["buffered_centroid"] = gdf["centroid"].buffer(2000)
```

→ Creating buffer around zip areas

→ Creating buffer around centroids

```
1 ax = gdf["buffered"].plot(alpha=.5)
2 gdf["buffered_centroid"].plot(ax=ax, color="red", alpha=.5)
3 gdf["boundary"].plot(ax=ax, color="white", linewidth=.5)
```

<AxesSubplot:>





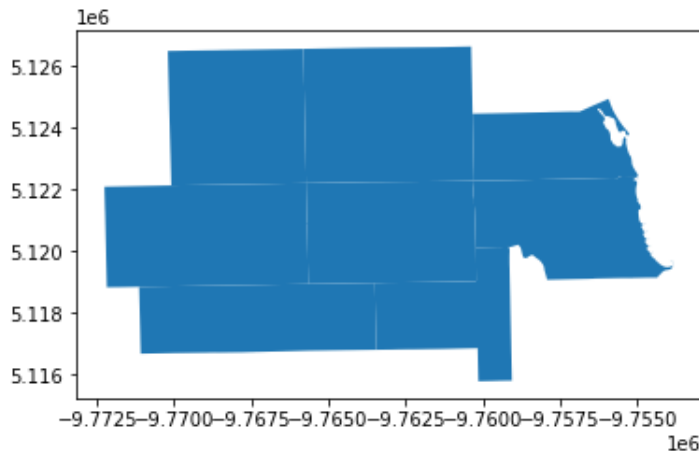
# Geometry relations

- GeoPandas offers a series of operations for geometry relations:
  - Crosses, intersects, overlaps, covers, within, touches, ...

```
1 selected = gdf.iloc[0]['geometry']
```

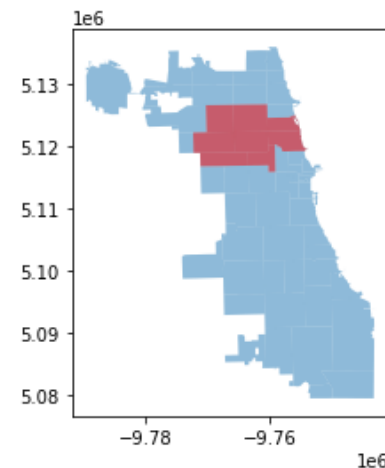
```
1 gdf[gdf['buffered'].intersects(selected)].plot()
```

<AxesSubplot:>



```
1 intersected = gdf[gdf['buffered'].intersects(selected)]
2 ax = gdf.plot(alpha=.5)
3 intersected.plot(ax=ax, color="red", alpha=.5)
```

<AxesSubplot:>



# Geometry relations

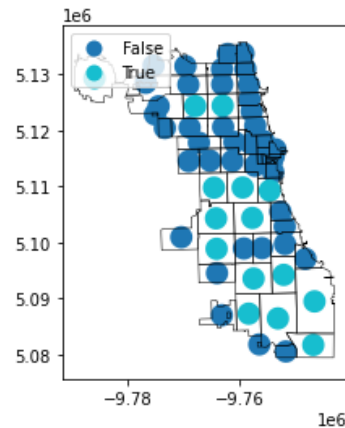
```
1 gdf["within"] = gdf["buffered_centroid"].within(gdf)
2 gdf["within"]
```

Whether buffered centroid is within zip area

```
0 False
1 False
2 False
3 False
4 False
...
56 True
57 True
58 True
59 False
60 True
Name: within, Length: 61, dtype: bool
```

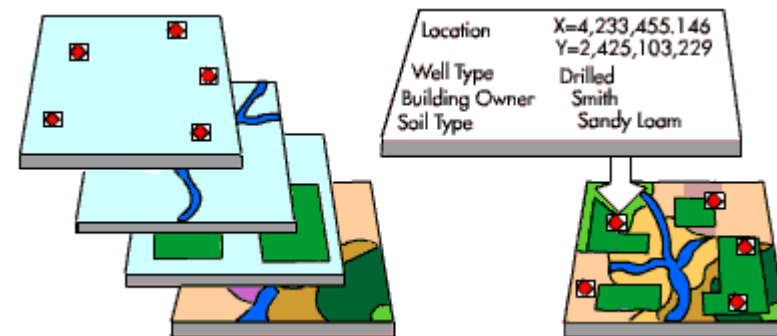
```
1 gdf = gdf.set_geometry("buffered_centroid")
2 ax = gdf.plot("within", legend=True, categorical=True, legend_kwds={'loc': "upper left"})
3 gdf["boundary"].plot(ax=ax, color="black", linewidth=.5)
```

<AxesSubplot:>



# Spatial joins

- A spatial join combines two GeoDataFrames based on the spatial relationship between their geometries.
- Example: spatial join between point layer (e.g., taxi pickups) and a polygon layer (e.g., zip codes).



# Spatial join

```
1 x_min, y_min, x_max, y_max = gdf.total_bounds
2
3 n = 1000
4
5 x = np.random.uniform(x_min, x_max, n)
6 y = np.random.uniform(y_min, y_max, n)
7
8 gdf_points = gpd.GeoDataFrame(geometry=gpd.points_from_xy(x, y), crs=gdf.crs)
9 gdf_points = gdf_points[gdf_points.within(gdf.unary_union)]
```

Bounds of GeoDataFrame

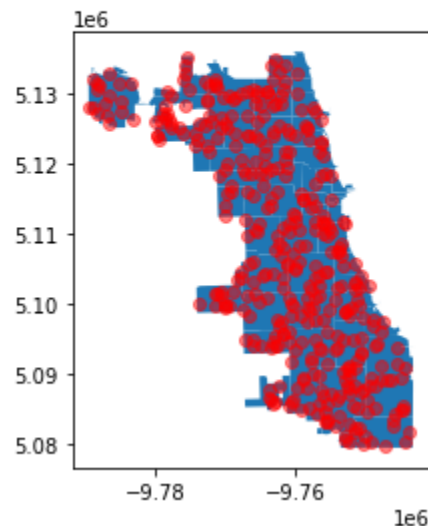
Sample size

Creating GeoDataFrame

Only keeping points within polygons

```
1 ax = gdf.plot()
2 gdf_points.plot(ax=ax, color="red", alpha=.5)
```

<AxesSubplot:>



# Spatial join

- Spatial join: for each point, is it *within* what zip code?

```
1 gpd.sjoin(gdf_points, gdf, predicate='within')
```

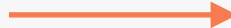
	geometry	index_right	objectid	shape_area	shape_len	zip
1	POINT (-9774274.365 5122584.288)	51	52	194062612.162	77647.3180069	60634
133	POINT (-9774348.810 5124145.541)	51	52	194062612.162	77647.3180069	60634
164	POINT (-9776525.134 5125184.977)	51	52	194062612.162	77647.3180069	60634
207	POINT (-9770521.404 5126432.046)	51	52	194062612.162	77647.3180069	60634
253	POINT (-9777824.534 5126280.441)	51	52	194062612.162	77647.3180069	60634
...	...	...	...	...	...	...
578	POINT (-9761501.347 5118140.579)	3	36	70853834.3797	42527.9896789	60622
917	POINT (-9763366.621 5118966.647)	3	36	70853834.3797	42527.9896789	60622
714	POINT (-9755003.744 5114813.246)	40	26	4847124.8171	14448.1749926	60602
756	POINT (-9772253.752 5121497.783)	2	35	45069038.4783	27288.6096123	60707
894	POINT (-9759652.063 5132688.902)	8	1	49170578.9623	33983.9133065	60626

392 rows × 6 columns

# Spatial join

- Grouping by zip code value to obtain number of points within that area.

```
1 result = gpd.sjoin(gdf_points, gdf, predicate='within').groupby('zip').count()
```

 Grouping by zip code

```
1 result = result.filter(['geometry'])  
2 result = result.rename(columns={'geometry': 'count'})
```

```
1 result
```

count	
zip	
60602	1
60605	1
60607	7
60608	12
60609	11
60610	3
60612	4
60613	4
60614	2
60615	8
60616	12
60617	20

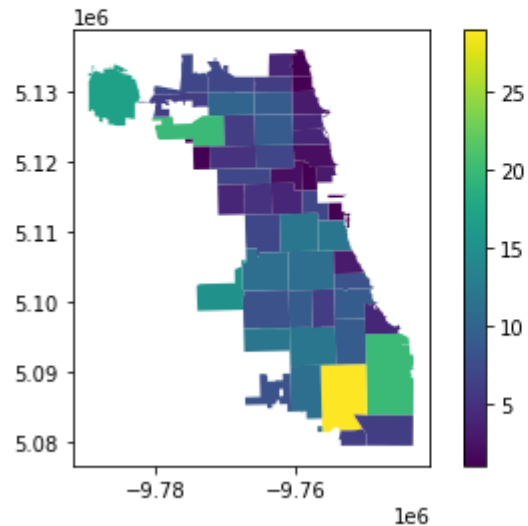
# Spatial join

```
1 merged = pd.merge(result, gdf, right_on='zip', left_index=True)
```

—————→ Merging with previous zip  
code GeoDataFrame

```
1 merged = merged.set_geometry('geometry')  
2 merged.plot('count', legend=True)
```

<AxesSubplot:>



# Spatial join

Aggregating points over spatial regions:

1. Spatial join: map between point and polygon.
2. Group by: aggregate (sum, count, mean, ...) by polygon.
3. Merge / join: map between aggregations and polygons.