

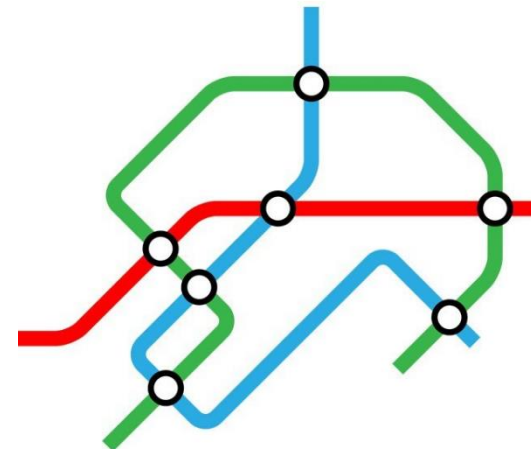
# Working with Spatial Data. Network Analysis

Reading, exploring and analyzing,  
feature extraction

**Yordan Darakchiev**

Technical Trainer

[iordan93@gmail.com](mailto:iordan93@gmail.com)





sli.do

#DataScience

# Table of Contents

- Geospatial data
  - Reading and exploring
  - Projections
  - Visualization
    - Scatter plots
    - Choropleth maps
- Network analysis
  - Graphs, types of graphs
  - Shortest path between nodes
  - Centrality
  - Communities

# Geospatial Data

Exploring, analyzing  
and visualizing

# Geospatial Data

- Data that has a geographic component to it
  - Most commonly: coordinates (latitude, longitude)
  - Sometimes: country, city, ZIP code, address
  - Not necessarily on Earth ([Google Mars](#))
- Sources
  - Satellite images
  - GPS data
  - Geotagging (e.g., photos at Facebook)
  - Manual entry, etc.
- Working with spatial data isn't trivial...
  - E.g., we need geometry on a sphere to calculate distances
  - ... but we have libraries that make our lives easier

# Reading and Exploring Geospatial Data

- In some cases, we have convenient datasets
- In other cases, it's in specific formats
  - GeoJSON, Shapefile, KML, etc.
  - Some libraries (like geopandas) can read these automatically
- Data cleaning
  - Non-spatial columns: proceed as usual
    - Tidy up the data, impute or remove missing values, explore outliers, normalize columns, etc.
  - Spatial columns: fixing or changing coordinates is easier when you visualize them
- Exploratory data analysis
  - Most commonly: look for clusters and other patterns
  - Also: compare attributes across different regions
    - E.g., income by country

# Example: Earthquake Data

- Dataset: earthquakes.csv, [info](#)
  - Read the dataset, look at missing values
  - Leave only columns you're interested in
- Explore the dataset
  - Examples: how is the magnitude distributed? When and where did the most powerful earthquakes happen? What are the recent ones?
- Perform additional data cleaning, exploration and visualization of the non-spatial columns
- Fix dates (remove invalid date format, convert to datetime)

```
dt_info = earthquake_data.Date + " " + earthquake_data.Time
earthquake_data = earthquake_data.drop(
    index = dt_info[dt_info.str.len() > 20].index)
earthquake_data["DateTime"] = pd.to_datetime(
    earthquake_data.Date + " " + earthquake_data.Time)
```

# Plotting Data on a Map

- To plot data, we'll use the basemap package

```
conda install -c anaconda basemap
```

```
from mpl_toolkits.basemap import Basemap
```

- Setting up and displaying a world map

```
m = Basemap(projection = "merc", llcrnrlat = -80, urcrnrlat = 80,  
            llcrnrlon = -180, urcrnrlon = 180)  
m.drawcoastlines()  
m.fillcontinents(color = "coral", lake_color = "aqua")  
m.drawparallels(np.arange(-90, 91, 30))  
m.drawmeridians(np.arange(-180, 181, 60))  
m.drawmapboundary(fill_color = "aqua")  
plt.show()
```

- Projections ([docs](#))

- Different ways to show a sphere in a 2D plane
- **Every projection has distortions**



# Plotting Data on a Map (2)

- Convert geographic coordinates ( $\varphi, \lambda$ ) to Cartesian coordinates ( $x, y$ )

- $x, y$  are measured in meters

```
x, y = m(earthquake_data.Longitude.tolist(),  
         earthquake_data.Latitude.tolist())
```

- Plot the coordinates ( $x, y$ ) on the map

```
m.plot(x, y, "o", markersize = 2, color = "red")
```

- Draw the other parts of the map

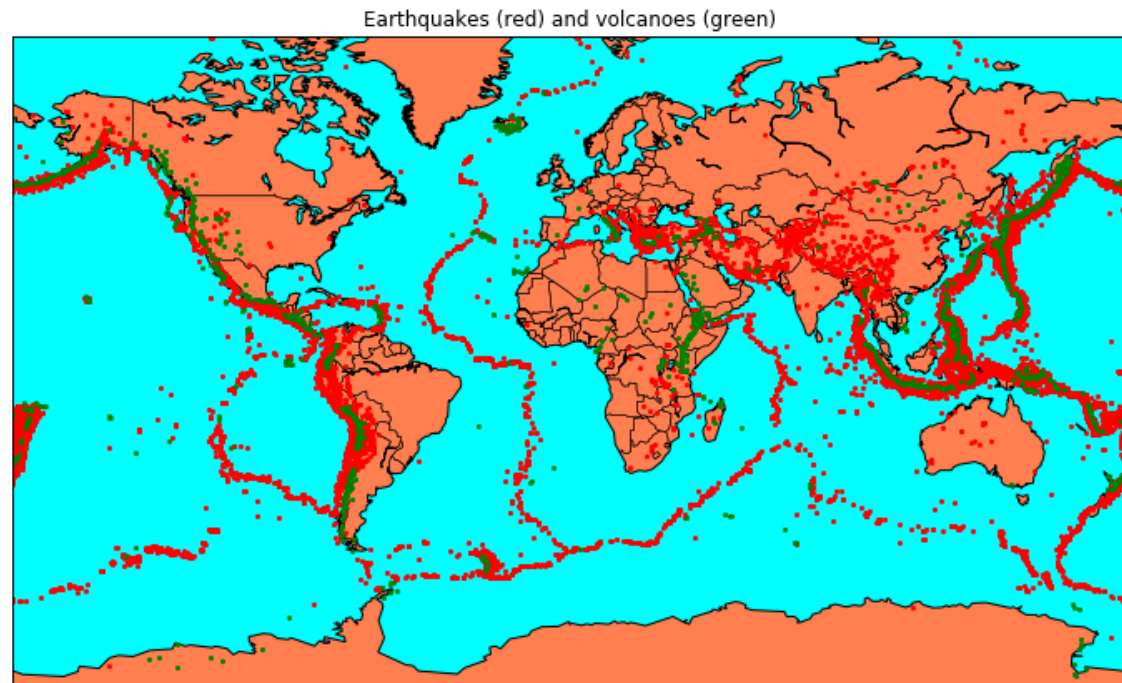
- Continents, countries, water

```
m.drawcoastlines()  
m.drawcountries()  
m.fillcontinents(color = "coral", lake_color = "aqua")  
m.drawmapboundary(fill_color = "aqua")  
m.drawcountries()  
plt.show()
```

# Adding Data on Volcanoes

- Dataset: `volcanoes.csv`, [info](#)
- Read the data and convert to x, y coordinates
- Plot just after the earthquakes
  - And before the "map decorations"

```
x_volc, y_volc = m(volcanos_data.Longitude.tolist(),  
                  volcanos_data.Latitude.tolist())  
m.plot(x_volc, y_volc, "o",  
       markersize = 4, color = "green")
```



# Drawing a Choropleth Map

- Like a heatmap
  - Shows different countries (or US states) in different colors according to a scale
- Dataset: `ufo_sightings_scrubbed.csv`, [info](#)
  - Clean the data (careful with "longitude")
  - Narrow down the data to US

```
ufos = pd.read_csv("ufo_sightings_scrubbed.csv", low_memory = False)
ufos = ufos[["datetime", "country", "state", "latitude", "longitude "]]
ufos.columns = ["datetime", "country", "state", "latitude", "longitude"]

ufos = ufos[ufos.country == "us"]
```

- Download the 3 shape files from [here](#) (st99\_00)

# Drawing a Choropleth Map (2)

- Create a map and read the shape file

```
m = Basemap(projection = "merc", llcrnrlon = -130, llcrnrlat = 23,  
            urcnrlon = -64, urcnrlat = 50)  
us_info = m.readshapefile("st99_d00", "states", drawbounds = True)
```

- Read the state names from `state_names.csv`
  - Use them to add the full names to the UFOs dataset

```
state_names = pd.read_csv("states.csv")  
state_names.abbreviation = state_names.abbreviation.str.lower()  
state_names_dict = {state.abbreviation: state["name"]  
                    for index, state in state_names.iterrows()}  
  
ufos.state.replace(state_names_dict, inplace = True)
```

- Get the number of sightings per state

```
num_sightings_by_state = ufos.groupby("state").size()
```

# Drawing a Choropleth Map (3)

- Import some libraries

```
import matplotlib
from matplotlib.colors import rgb2hex
from matplotlib.patches import Polygon
```

- Set up the map and some objects to use later

```
fig = plt.figure(figsize = (15, 10))
m = Basemap(projection = "merc", llcrnrlon = -130, llcrnrlat = 23,
            urcrnrlon = -64, urcrnrlat = 50)
us_info = m.readshapefile("st99_d00", "states", drawbounds = True)
colors = {}
state_names = []
cmap = plt.cm.Greens
vmin = num_sightings_by_state.min()
vmax = num_sightings_by_state.max()
```

# Drawing a Choropleth Map (4)

- Compute colors for each state
  - Using a specified color map
  - `np.sqrt()` spreads the colors more evenly
  - `(sightings - vmin) / (vmax - vmin)` returns a normalized value from 0 to 1
  - `cmap()` returns RGBA values, `[:3]` discards the alpha channel

```
for shape_dict in m.states_info:
    state_name = shape_dict["NAME"]
    # Skip DC and Puerto Rico
    if state_name not in ["District of Columbia", "Puerto Rico"]:
        sightings = num_sightings_by_state[
            num_sightings_by_state.index == state_name][0]
        colors[state_name] = cmap(
            np.sqrt((sightings - vmin) / (vmax - vmin)))[:3]
        state_names.append(state_name)
```

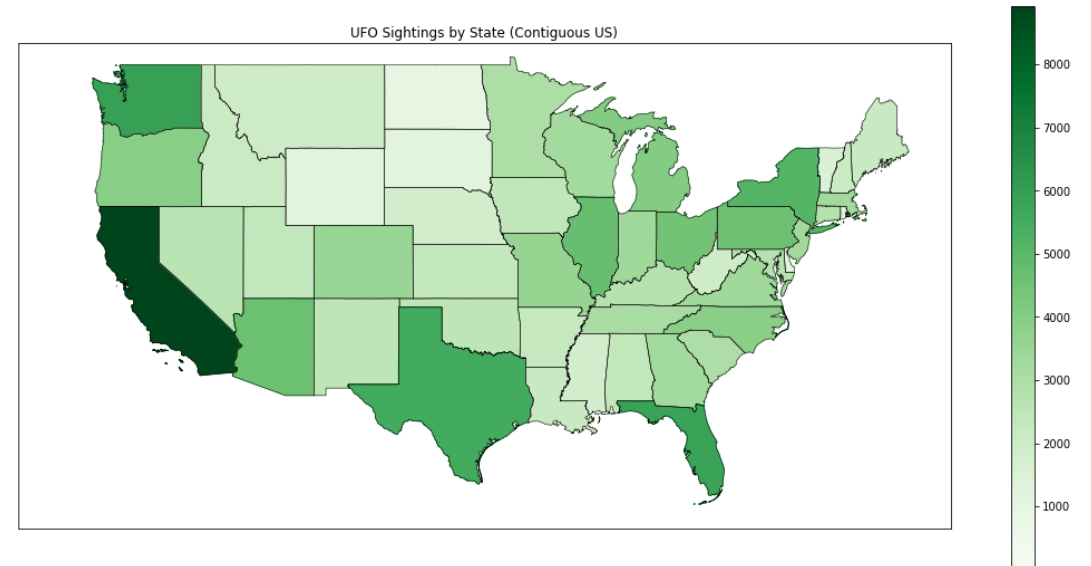
# Drawing a Choropleth Map (5)

- Draw the polygons for each state

```
ax = plt.gca()
for nshape, seg in enumerate(m.states):
    # Skip DC and Puerto Rico
    if state_names[nshape] not in ["District of Columbia", "Puerto Rico"]:
        color = rgb2hex(colors[state_names[nshape]])
        poly = Polygon(seg, facecolor = color, edgecolor = color)
        ax.add_patch(poly)
```

- Add title and color bar

```
plt.title(
    "UFO Sightings by State (Contiguous US)")
colorbar_ax = fig.add_axes(
    [0.95, 0.15, 0.02, 0.7])
matplotlib.colorbar.ColorbarBase(
    colorbar_ax, cmap = cmap,
    norm = matplotlib.colors.Normalize(
        vmin, vmax))
plt.show()
```



# Analyzing Maps

- There are many algorithms used to model spatial data
  - Most commonly, we look for density patterns and clusters of points
  - Common algorithms are
    - [KDE](#) – Kernel Density Estimation
    - [kMeans](#) Clustering
    - [Hierarchical](#) Clustering
    - [kNN](#) – k Nearest Neighbors
  - This course doesn't deal with modelling, so we won't get into more detail
    - But feel free to explore the algorithms as you wish
    - You can see details on these on machine learning-related articles
- We can also represent the map as a network
  - E.g. road maps, railway maps, or other "sets of connected dots"



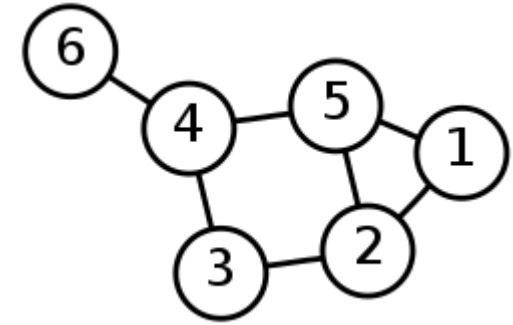


# Network Analysis

## Working with graphs

# Networks = Graphs

- A graph is a geometrical object consisting of objects which are related by some attribute
  - **Nodes** (vertices, points) – describe objects
  - **Edges** (arcs, lines) – connect nodes

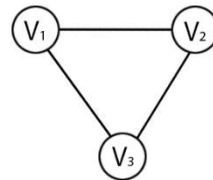


- Types of graphs

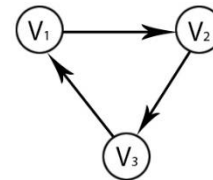
- Directed / undirected

- In a directed graph, there is only one way to travel between the nodes

Undirected Graph

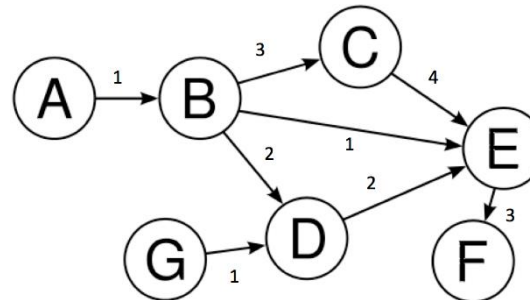


Directed Graph




- Weighted / unweighted

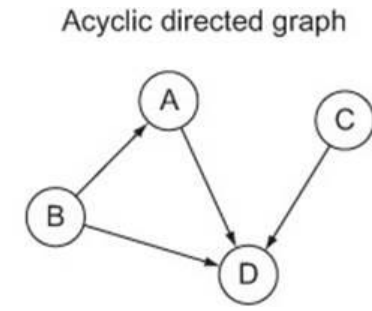
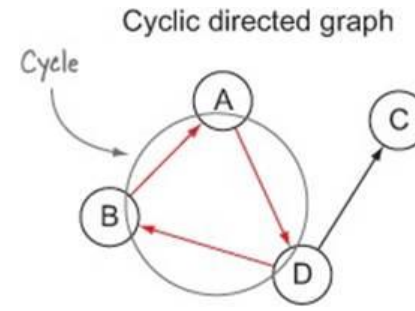
- A weighted graph contains some quantity ("weight", usually  $\geq 0$ ) over each of its edges



# Graphs

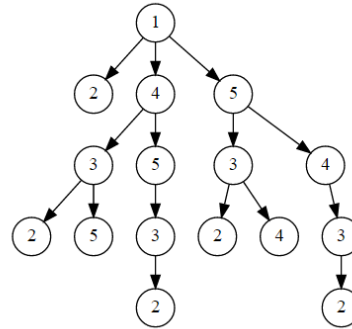
## ■ Types of graphs (cont'd)

- **Cyclic / acyclic**
    - When you travel along a cyclic graph, you will visit one node more than once
  - These types are independent
    - i.e. a graph can be “acyclic directed unweighted graph”
- 
- ```
graph LR; A((A)) -- red --> B((B)); B -- red --> D((D)); D -- red --> A; D -- black --> C((C));
```

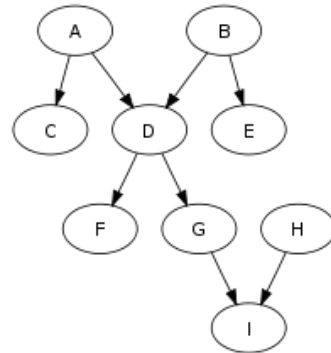


- Special cases

- **Tree** – each node has at most one "parent"



- **DAG** – directed acyclic graph



# Representing Graphs

- We can use the library networkx
  - Installed by default with Anaconda
- Create a simple weighted undirected graph

```
import networkx as nx
g = nx.Graph()
g.add_edge("a", "b", weight = 0.1)
g.add_edge("b", "c", weight = 1.5)
g.add_edge("a", "c", weight = 1.0)
g.add_edge("c", "d", weight = 2.2)
```

- Display the graph

```
nx.draw(g, with_labels = True)
plt.show()
```

# Finding a Shortest Path

- Advanced graph display
  - Show the weights at each edge
  - Make the edge width proportional to its weight

```
pos = nx.spring_layout(g)
weights = nx.get_edge_attributes(g, "weight")
nx.draw(g, pos, with_labels = True)

nx.draw_networkx_edge_labels(g, pos,
    edge_labels = weights)
nx.draw_networkx_edges(g, pos,
    width = [v * 2 for v in weights.values()])
plt.show()
```

- Shortest paths

```
print(nx.shortest_path(g, "b", "d"),
    nx.shortest_path_length(g, "b", "d"))
print(nx.shortest_path(g, "b", "d", weight = "weight"),
    nx.shortest_path_length(g, "b", "d", weight = "weight"))
```

# Creating Directed Graphs

- Directed graph (digraph)
  - Simply change the definition of g
  - Now each edge is directed
  - The visualization will include arrows
    - They point at the direction of each connection

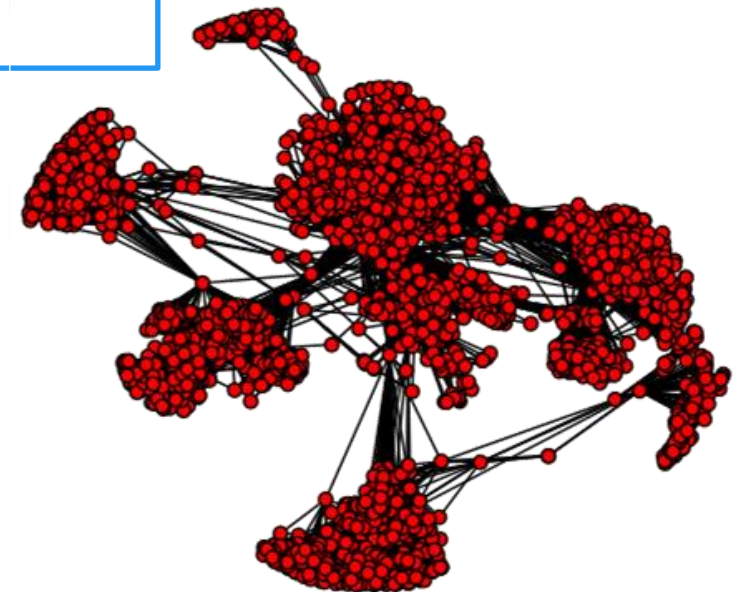
```
g = nx.DiGraph()
g.add_edge("a", "b", weight = 0.1)
g.add_edge("b", "c", weight = 1.5)
g.add_edge("a", "c", weight = 1.0)
g.add_edge("c", "d", weight = 2.2)
```

```
print(nx.shortest_path(g, "b", "d")) # ['b', 'c', 'd']
print(nx.shortest_path(g, "d", "b")) # Error: No path between d and b.
```

# Example: Social Circles

- Dataset: facebook.zip, [info](#)
  - Format: first\_user\_id second\_user\_id
    - I.e. edge list
- Read the graph
  - Extremely simple

```
facebook_graph = nx.read_edgelist("facebook_combined.txt")  
print(len(facebook_graph.nodes)) # 4039  
print(len(facebook_graph.edges)) # 88234
```

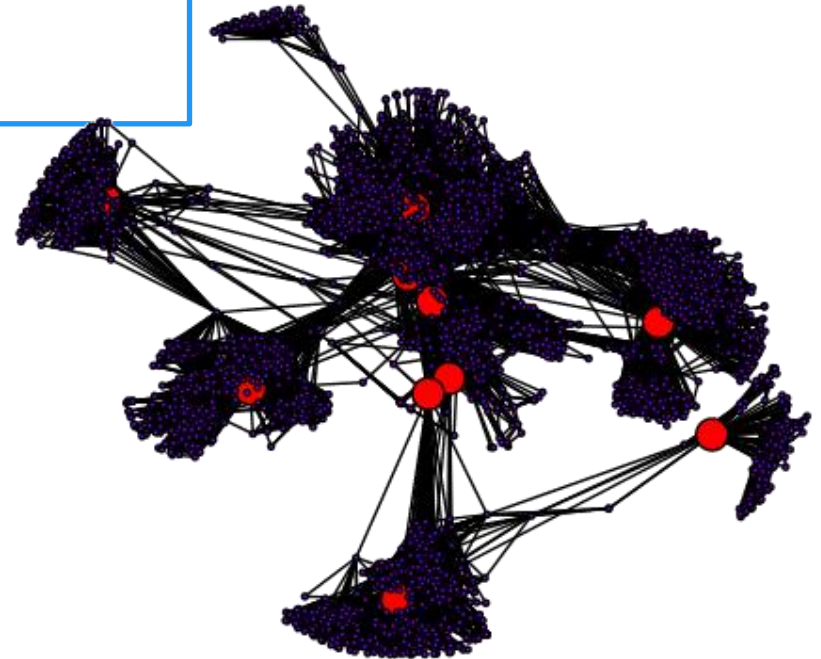


# Calculating Important Nodes

- Measure: **centrality**
  - Different types of centrality, according to different formulas
    - E.g. "betweenness centrality"
    - Measures how important a node is
- To exemplify, let's use a smaller graph

```
karate_graph = nx.karate_club_graph()  
centrality = nx.betweenness_centrality(karate_graph)  
# Returns a dictionary
```

- Ten most important nodes in the Facebook graph
  - Look similar to cluster centroids



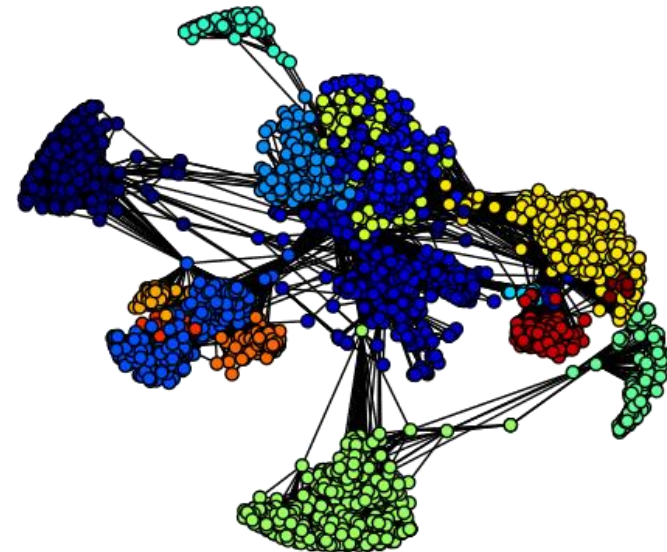


# Finding Communities

- Measure: **cliques**
  - Most commonly used algorithm: Girvan – Newman
    - Uses **edge betweenness** as the measure

```
from networkx.algorithms import community
nx.draw(karate_graph, with_labels = True)
communities_generator = community.girvan_newman(karate_graph)
for i in range(1, 4):
    communities = next(communities_generator)
    print("level " + str(i), communities)
```

- We can find communities in the Facebook graph
  - Look similar to different clusters



# Summary

- Geospatial data
  - Reading and exploring
  - Projections
  - Visualization
    - Scatter plots
    - Choropleth maps
- Network analysis
  - Graphs, types of graphs
  - Shortest path between nodes
  - Centrality
  - Communities

The image features a white background with two blue decorative bars. The top bar is a solid blue strip. The bottom bar is a gradient blue strip that transitions from a lighter blue on the left to a darker blue on the right. The word "Questions?" is centered in a blue, sans-serif font.

Questions?