

Python for Geospatial Data Analysis

CoPhil EO AI/ML Training - Day 1, Session 3

Stylianos Kotsopoulos
EU-Philippines CoPhil Programme

Welcome to Session 3

Session Overview

Python for Geospatial Data Analysis

Hands-on introduction to working with vector and raster data using Python

Format: Brief conceptual intro + Extended hands-on coding


Duration: 2 hours (15-20 min presentation + 100 min hands-on)

Learning Objectives

By the end of this session, you will be able to:

1. Set up and use Google Colaboratory for geospatial analysis
2. Load, explore, and visualize vector data with **GeoPandas**
3. Read, process, and visualize raster data with **Rasterio**
4. Perform basic geospatial operations (clipping, reprojecting, cropping)
5. Prepare data for AI/ML workflows

Session Roadmap

Time	Topic	Duration
00-15 min	Setup & Python Basics Recap	15 min
15-55 min	GeoPandas for Vector Data (HANDS-ON)	40 min
55-60 min	 Break	5 min
60-110 min	Rasterio for Raster Data (HANDS-ON)	50 min
110-120 min	Summary & Next Steps	10 min

Notebook Access

 Google Colab Notebook:

[Day1_Session3_Python_Geospatial_Data.ipynb](#)

1. **Open link** from course materials
2. **Save a copy** to your Drive
3. **Run first cell** to install packages
4. **Follow along** as we code together

Today's Focus

Vector Data:

- Administrative boundaries
- Points of interest
- Roads, rivers
- Training sample polygons
- Using **GeoPandas**

Raster Data:

- Satellite imagery
- Digital elevation models
- Land cover maps
- AI model outputs
- Using **Rasterio**

Integration: Combining vector and raster for complete EO workflows

Why Python for Geospatial?

The Python Advantage

Why Python is the Leading Language for EO:

1. Rich Ecosystem

- Hundreds of specialized libraries
- Active development and community

2. Easy to Learn

- Clear syntax, readable code
- Gentle learning curve

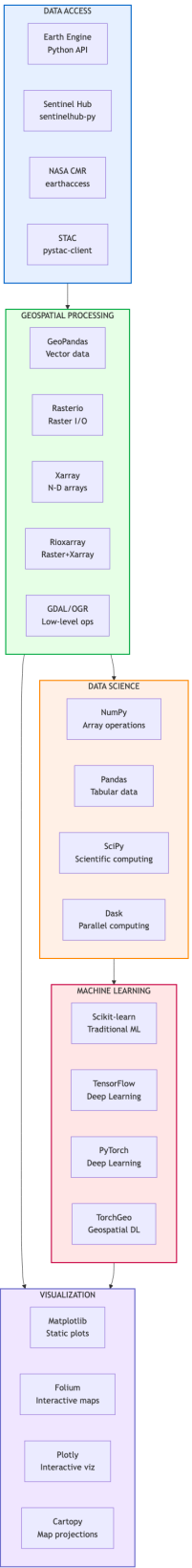
3. Powerful Integration

- Connects data sources, processing, ML
- Single environment for complete workflows

4. Free and Open Source

- No licensing costs
- Transparent and reproducible

Python Geospatial Ecosystem



Complete Python Earth Observation Ecosystem organized by function

Integration Capabilities

Python Connects Everything:

```
1 # Example workflow
2 import geopandas as gpd
3 import rasterio
4 from sklearn.ensemble import RandomForestClassifier
5
6 # Load vector training data
7 training = gpd.read_file('samples.geojson')
8
9 # Load satellite raster
10 with rasterio.open('sentinel2.tif') as src:
11     image = src.read()
12
13 # Extract features and train model
14 X, y = extract_features(image, training)
15 model = RandomForestClassifier()
16 model.fit(X, y)
17
18 # Predict on full image
```

Community and Resources

Vibrant Python Geospatial Community:

Documentation:

- Comprehensive guides for all libraries
- Tutorials and examples
- API references

Community Support:

- Stack Overflow
- GitHub discussions
- GIS Stack Exchange
- Dedicated forums

Learning Resources:

- Free courses (Coursera, Udemy)
- Books (Automating GIS Processes)
- Blogs and tutorials

Conferences/workshops

Google Colaboratory

Why Colab for This Training?

Advantages for Learning:

1. No Setup Hassles

- Works immediately
- No environment configuration
- Consistent for all participants

2. Accessible Anywhere

- Just need a browser
- Works on any computer
- Even tablets

3. Powerful Resources

- Free GPU for deep learning
- 12+ GB RAM
- Sufficient for all exercises

4. Easy Sharing

- Share notebooks instantly

Collaborative editing

Colab Interface Overview

Main Components:

1. Menu Bar

- File, Edit, View, Insert, Runtime

2. Toolbar

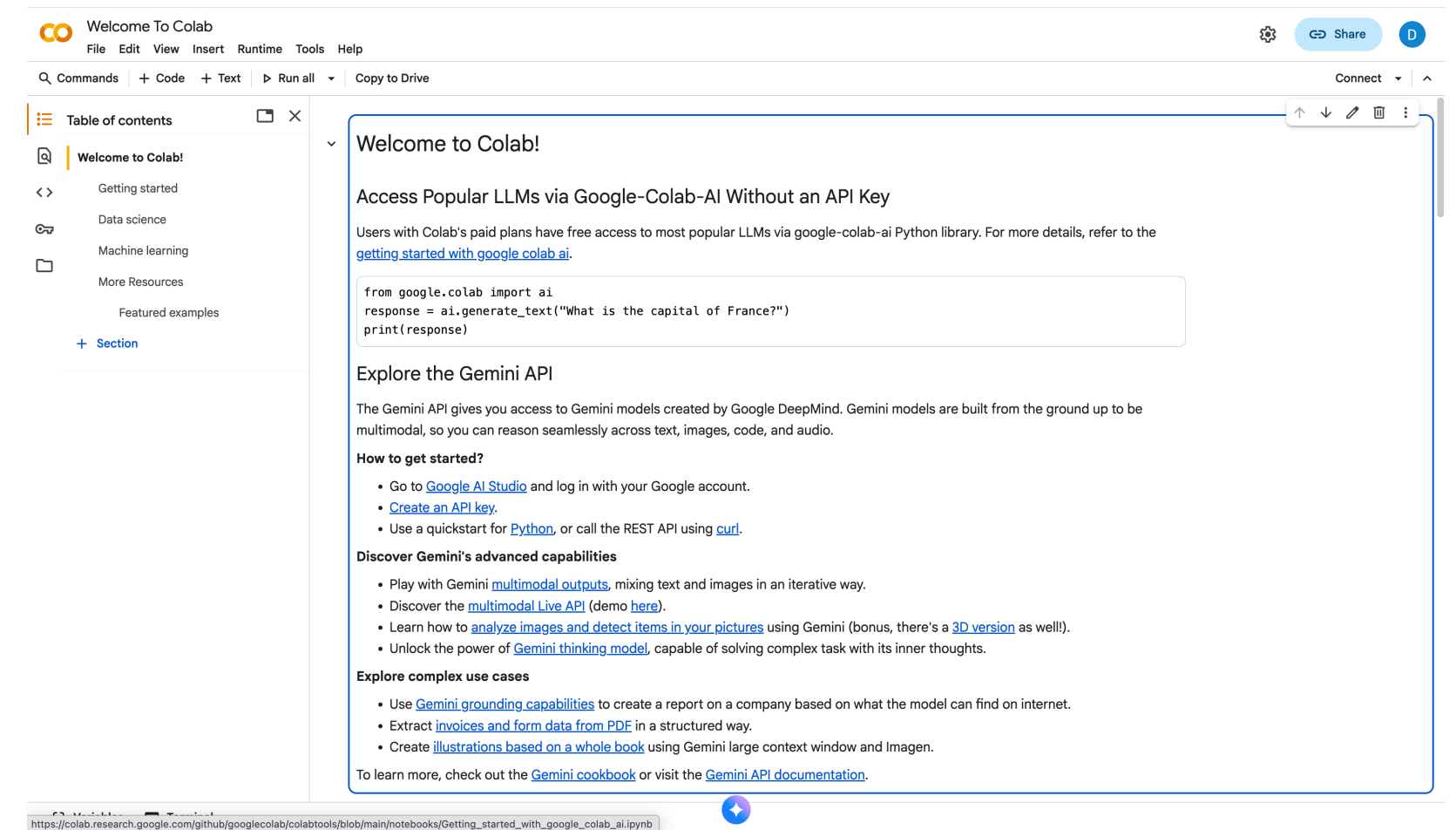
- Play button to run cells
- Add code/text cells

3. Notebook Area

- Code cells (executable)
- Text cells (Markdown)

4. Sidebar

- Table of contents
- Files browser
- Code snippets



Colab Interface

Running Code in Colab

Two Ways to Execute Cells:

1. Click the Play Button

- Left side of each code cell
- Runs that specific cell

2. Keyboard Shortcuts

- **Shift + Enter:** Run cell and move to next
- **Ctrl + Enter:** Run cell, stay on current
- **Ctrl + M then A:** Add cell above
- **Ctrl + M then B:** Add cell below

Output Appears Below Cell:

Text, plots, tables, errors all display inline.

Google Drive Integration

Mounting Your Google Drive:

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Benefits:

- Persistent storage (Colab resets)
- Upload/download data
- Save outputs
- Share files between notebooks

Access Your Files:

```
1 # Your Drive files appear at:
2 # /content/drive/MyDrive/
3
4 # Example:
5 import geopandas as gpd
6 gdf = gpd.read_file('/content/drive/MyDrive/data/boundaries.shp')
```

Installing Additional Packages

Most Common Libraries Pre-Installed:

NumPy, Pandas, Matplotlib, Scikit-learn

For Geospatial Libraries:

```
1 # GeoPandas (usually pre-installed, but check version)
2 !pip install geopandas
3
4 # Rasterio
5 !pip install rasterio
6
7 # Other useful libraries
8 !pip install earthengine-api
9 !pip install folium
```

Note: Packages need reinstalling each session (Colab resets runtime)

GeoPandas for Vector Data

What is GeoPandas?

Pandas + Geometry = GeoPandas

Definition:

Extension of Pandas for working with geospatial vector data

Key Concept:

Like a spreadsheet/table where one column contains geometries (points, lines, polygons)

Built On:

- **Pandas** - Data manipulation
- **Shapely** - Geometric operations
- **Fiona** - File I/O
- **PyProj** - Coordinate systems

GeoPandas

Pandas + Geometry

The GeoDataFrame Concept

Similar to Pandas DataFrame:

Regular DataFrame:

Name	Population	Area
Manila	1.78M	42.88
Cebu	0.92M	315
Davao	1.63M	2444

GeoDataFrame:

Name	Population	Area	geometry
Manila	1.78M	42.88	POLYGON(...)
Cebu	0.92M	315	POLYGON(...)
Davao	1.63M	2444	POLYGON(...)

Special “geometry” Column:

Contains spatial information (coordinates defining shapes)

Common Vector Data Operations

What You Can Do with GeoPandas:

1. Read/Write

- Shapefiles, GeoJSON, GeoPackage, PostGIS

2. Explore

- View attributes, examine geometries

3. Visualize

- Quick plotting, interactive maps

4. Geoprocessing

- Buffer, intersection, union, clip

5. Spatial Joins

- Combine datasets based on location

6. Coordinate Transformations

- Reproject to different CRS

GeoPandas Code Example

Loading and Visualizing:

```
1 import geopandas as gpd
2 import matplotlib.pyplot as plt
3
4 # Read Philippine provinces shapefile
5 provinces = gpd.read_file('philippines_provinces.shp')
6
7 # Examine data
8 print(provinces.head())
9 print(provinces.crs) # Check coordinate system
10
11 # Simple plot
12 provinces.plot(figsize=(10, 10),
13                 edgecolor='black',
14                 facecolor='lightblue')
15 plt.title('Philippine Provinces')
16 plt.show()
```

Visualization Capabilities

GeoPandas Plotting:

```
1 # Color by attribute
2 provinces.plot(column='population',
3               cmap='YlOrRd',
4               legend=True,
5               figsize=(12, 10))
6 plt.title('Population by Province')
7
8 # Add basemap (with contextily)
9 import contextily as ctx
10 provinces_web_mercator = provinces.to_crs(epsg=3857)
11 ax = provinces_web_mercator.plot(figsize=(15, 15),
12                               alpha=0.5)
13 ctx.add_basemap(ax)
```


Philippine Coordinate Reference Systems

Common CRS for Philippines:

EPSG	Name	Type	Units	Use Case
4326	WGS84	Geographic	Degrees	Web maps, lat/lon
32651	UTM Zone 51N	Projected	Meters	Western PH, Manila
32652	UTM Zone 52N	Projected	Meters	Eastern PH, Mindanao
3123	PRS92 Zone III	Projected	Meters	National standard

Rule: Use **geographic (4326)** for storage, **projected (UTM)** for analysis

Philippine UTM Zones

UTM Zone 51N (EPSG:32651)

Coverage: - Metro Manila - Palawan - Western Luzon - Western Visayas

Most common for: - Urban analysis - Palawan studies - Manila projects

Code Example:

```
1 # Reproject to UTM 51N for area calculation
2 gdf_utm = gdf.to_crs(epsg=32651)
3 gdf_utm['area_km2'] = gdf_utm.geometry.area / 1_000_000
```

UTM Zone 52N (EPSG:32652)

Coverage: - Mindanao - Eastern Visayas - Bicol Region - Eastern Luzon

Most common for: - Mindanao analysis - Disaster mapping - Agricultural studies

Philippine Geospatial Data Sources

NAMRIA Geoportal - Administrative boundaries -
Topographic maps - Land cover 2020 - DEMs -
www.geoportal.gov.ph

PhilSA - Satellite imagery - EO products - philsa.gov.ph

All work with GeoPandas - just `gpd.read_file()`!

PSA - Census boundaries - Barangay data - psa.gov.ph

OpenStreetMap - Roads, buildings - extract.bbbike.org

Natural Earth - Country boundaries -
naturalearthdata.com

Rasterio for Raster Data

What is Rasterio?

Python Wrapper for GDAL

Definition:

Clean, idiomatic Python library for reading and writing geospatial raster data

Why Not Use GDAL Directly?

- GDAL Python bindings are cumbersome
- Rasterio is more Pythonic
- Better integration with NumPy
- Cleaner syntax

Works With:

All formats GDAL supports - GeoTIFF, COG, NetCDF, HDF, etc.

Rasterio

Python Wrapper for GDAL

Raster Data Structure

How Rasterio Represents Imagery:

3D NumPy Array:

(bands, rows, columns)

Example: Sentinel-2 10m Bands:

```
1 # 4 bands (Blue, Green, Red, NIR)
2 # 1098 rows (10980 m / 10 m)
3 # 1098 columns (10980 m / 10 m)
4 # Shape: (4, 1098, 1098)
5
6 array[0, :, :] # Band 1 (Blue)
7 array[1, :, :] # Band 2 (Green)
8 array[2, :, :] # Band 3 (Red)
9 array[3, :, :] # Band 4 (NIR)
```

Reading Raster Data with Rasterio

Basic Workflow:

```
1 import rasterio
2 import numpy as np
3
4 # Open raster file
5 with rasterio.open('sentinel2_10m.tif') as src:
6     # Read all bands
7     data = src.read()
8
9     # Read specific band
10    red_band = src.read(3)
11
12    # Get metadata
13    print(f"CRS: {src.crs}")
14    print(f"Transform: {src.transform}")
15    print(f"Width: {src.width}, Height: {src.height}")
16    print(f"Bounds: {src.bounds}")
17    print(f"Number of bands: {src.count}")
```

Array Operations

Rasterio + NumPy = Powerful Processing

```
1 # Calculate NDVI
2 with rasterio.open('sentinel2_10m.tif') as src:
3     red = src.read(3).astype(float)
4     nir = src.read(4).astype(float)
5
6 # NDVI formula
7 ndvi = (nir - red) / (nir + red + 1e-8) # Small value prevents division by zero
8
9 # Apply threshold
10 vegetation_mask = ndvi > 0.3
11
12 # Calculate statistics
13 print(f"Mean NDVI: {np.mean(ndvi):.3f}")
14 print(f"Vegetation pixels: {np.sum(vegetation_mask)}")
```


Common Spectral Indices

Key Indices for EO Analysis:

Index	Formula	Purpose	Range
NDVI	$(\text{NIR} - \text{Red}) / (\text{NIR} + \text{Red})$	Vegetation health	-1 to +1
EVI	$2.5 \times (\text{NIR} - \text{Red}) / (\text{NIR} + 6 \times \text{Red} - 7.5 \times \text{Blue} + 1)$	Enhanced vegetation	-1 to +1
NDWI	$(\text{Green} - \text{NIR}) / (\text{Green} + \text{NIR})$	Water bodies	-1 to +1
NDBI	$(\text{SWIR} - \text{NIR}) / (\text{SWIR} + \text{NIR})$	Built-up areas	-1 to +1

Sentinel-2 Bands: Blue (B2), Green (B3), Red (B4), NIR (B8), SWIR (B11, B12)

Philippine Application: Rice Monitoring

Using NDVI for Philippine Rice Paddies:

```
1 # Calculate NDVI for Central Luzon rice area
2 with rasterio.open('sentinel2_central_luzon.tif') as src:
3     red = src.read(4).astype(float)    # Band 4
4     nir = src.read(8).astype(float)    # Band 8
5
6 # NDVI calculation
7 ndvi = (nir - red) / (nir + red + 1e-8)
8
9 # Classify vegetation health
10 bare_soil = ndvi < 0.2    # Recently planted / fallow
11 growing = (ndvi >= 0.2) & (ndvi < 0.5) # Early growth
12 mature = (ndvi >= 0.5) & (ndvi < 0.8)  # Peak biomass
13 very_dense = ndvi >= 0.8    # Maximum vegetation
14
15 # Calculate rice area statistics
16 pixel_area = 100    # 10m x 10m = 100 m²
17 mature_rice_area_ha = np.sum(mature) * pixel_area / 10000
18
```

Visualization with Rasterio

Displaying Satellite Imagery:

```
1 import matplotlib.pyplot as plt
2 from rasterio.plot import show
3
4 # Open and display
5 with rasterio.open('sentinel2_10m.tif') as src:
6     # Show true color composite (RGB)
7     show((src, [3, 2, 1]), title='True Color')
8
9     # Show false color composite (NIR, Red, Green)
10    show((src, [4, 3, 2]), title='False Color (NIR-R-G)')
11
12 # Or read and plot with matplotlib
13 with rasterio.open('sentinel2_10m.tif') as src:
14     data = src.read([3, 2, 1]) # RGB
15     # Scale to 0-255 for display
16     data_scaled = np.clip(data / 3000, 0, 1)
17
18     plt.figure(figsize=(10, 10))
```

Hands-On Preview

What We'll Build Today

Notebook 1: Vector Data with GeoPandas

1. Load Philippine administrative boundaries
2. Explore and visualize provinces
3. Filter specific regions (e.g., Central Luzon)
4. Calculate area and basic statistics
5. Spatial operations (buffer, clip)
6. Create training sample polygons

What We'll Build Today

Notebook 2: Raster Data with Rasterio

1. Load Sentinel-2 image subset
2. Examine metadata and properties
3. Visualize true and false color composites
4. Calculate vegetation indices (NDVI, EVI)
5. Crop to area of interest
6. Extract pixel values at point locations
7. Save processed outputs

Integrating Vector and Raster

Combining Both Data Types:

```
1 import geopandas as gpd
2 import rasterio
3 from rasterio.mask import mask
4
5 # Load vector boundary
6 aoi = gpd.read_file('study_area.geojson')
7
8 # Load raster
9 with rasterio.open('sentinel2.tif') as src:
10     # Clip raster to vector boundary
11     clipped_data, clipped_transform = mask(
12         src,
13         aoi.geometry,
14         crop=True
15     )
16
17 # Update metadata for output
18 out_meta = src.meta.copy()
```

Preparing for ML Workflows

What You'll Learn:

1. Extract Training Data

- Sample raster values at polygon locations
- Create feature matrix (X) and labels (y)

2. Spatial Data Splits

- Avoid spatial autocorrelation in train/test

3. Data Formatting

- Structure for Scikit-learn, TensorFlow, PyTorch

4. Quality Control

- Check for NaN values, outliers
- Validate spatial alignment

Link to Colab Notebooks

Access Today's Notebooks:

Notebook 1: Vector Data

[Link will be provided in chat]

Notebook 2: Raster Data

[Link will be provided in chat]

Make a Copy:

File → Save a copy in Drive (so you can edit and experiment)

Tips for Success

As We Work Through Notebooks:

1. Run Cells Sequentially

- Top to bottom order matters
- Each cell may depend on previous

2. Read the Comments

- Explanations included in code
- Learn the “why” not just “how”

3. Experiment

- Modify parameters
- Try different visualizations
- Break things and learn

4. Ask Questions

- Use chat or raise hand
- No question is too basic

5. Take Notes

Useful patterns and code snippets

Key Concepts Summary

Python Geospatial Ecosystem

GeoPandas:

- DataFrame + geometry column
- Vector data operations
- Easy visualization
- Integration with Pandas

Rasterio:

- Clean GDAL wrapper
- NumPy array representation
- Comprehensive metadata handling
- Efficient I/O

Integration:

- Both work together seamlessly
- Complete EO workflows in Python
- Foundation for ML/AI applications

Why These Skills Matter

For AI/ML in Earth Observation:

1. Data Preparation

- Loading and preprocessing is 80% of work
- Quality in → Quality out

2. Feature Engineering

- Calculate indices, textures, derivatives
- NumPy operations on raster arrays

3. Training Data Creation

- Sample raster at polygon locations
- Extract features for supervised learning

4. Model Deployment

- Apply trained models to new imagery
- Generate prediction maps

5. Validation and QA

- Compare predictions to ground truth

Calculate accuracy metrics

Building Blocks for This Week

Today's Skills Enable:

Day 2:

- Random Forest land cover classification
- Preparing training data for ML

Day 3:

- Deep learning data pipelines
- U-Net flood mapping inputs

Day 4:

- Time series data preparation
- LSTM input formatting

Mastering these fundamentals now will make everything else smoother.

 5-Minute Break

Stretch Break

Stand up • Grab water • Back in 5 minutes

Let's Begin!

Transition to Hands-On

Open Your Notebooks

We'll start with:

Vector Data Analysis using GeoPandas

Remember:

- Make a copy of the notebook
- Mount your Google Drive
- Run cells in order
- Ask questions anytime

Support During Hands-On

Instructors Available:

- Main instructor demonstrating
- Teaching assistants in chat
- Screen sharing for debugging

Pacing:

- We'll work through together
- Pause points for questions
- Extra exercises for fast finishers

Goal:

Everyone completes core exercises, understands concepts, ready for GEE

Session Complete!

Session Summary

What You've Learned:

- ✓ Google Colab setup for geospatial work
- ✓ GeoPandas for vector data (load, visualize, analyze)
- ✓ Rasterio for raster data (read, process, visualize)
- ✓ Coordinate systems and projections
- ✓ Basic geospatial operations (clip, reproject, crop)
- ✓ Data preparation for ML workflows

Q&A

Common Questions:

- GeoPandas vs Shapely?
- When to use Rasterio vs GDAL?
- CRS issues and solutions?
- Memory errors with large files?
- Best practices for file paths?
- NoData values handling?
- Visualization tips?
- Integration with ML pipelines?

Next: Session 4

Google Earth Engine

Coming up:

- Cloud-based EO data processing
- Access to entire Sentinel archive
- Planetary-scale analysis
- Python API (geemap)
- Cloud masking & compositing
- Export workflows

Get ready for GEE! 🌍

Everyone completes core exercises with understanding, not just copying code.

Questions Before We Start?

Any questions about the tools or approach?

Let's Code!

Opening Notebook 1: GeoPandas for Vector Data