

Assignment 03: Rasters and Tables

For this assignment, you will be using Numpy, GDAL, and Pandas. Most of these questions can be solved using methods shown in lecture. However, some cases will require you to search the internet for answers. This is intended because efficiently searching documentation or Stackoverflow is a requirement of modern programming. This assignment also requires you to obtain data (I will tell you where) and organize it correctly so that you can access it via relative paths (do not use absolute paths anywhere in this assignment).

To turn this in, create a PDF of your finished Jupyter notebook using Quarto (we will discuss this on Tuesday) and upload it to Canvas.

Question 1

For this question, you will be using numpy to manipulate arrays in a way that works well for big data and will be using spatial tools to assess land-use change. Specifically, you will use data from a global, 300m resolution land-use, land-cover (LULC) map produced by ESACCI. They provide a wonderful time series of data from 1992-2019. If you want more information, check out their web viewer at <https://maps.elie.ucl.ac.be/CCI/viewer/>. Here you can view the data but also learn more about the classification regressions that uses time-series spectral data patterns to classify grid-cells into different types. I have extracted lulc maps for the country of Rwanda, which you will find in this assignment's data folder.

Part a.

Using the gdal package and the `gdal.Open()` function, open up the land-use, land-cover map for Rwanda in 2000.

In this file, there is only 1 band in this file, so you can also access it with the `GetRasterBand(1)` function. Without reading the whole array, show how you can determine how many total grid-cells there are in this country.

```

# 1A Answer
from osgeo import gdal
import numpy as np
import os
import matplotlib.pyplot as plt

# Set up the data directory path (relative to A3 directory)
# From A3: ../../../../base_data (up to Apec8222, then into base_data)
data_directory = '../../../../base_data'
lulc_filename_2000 = 'rwanda_lulc_2000.tif'
lulc_file_path_2000 = os.path.join(data_directory, lulc_filename_2000)

# Check if file exists
print(f'Looking for file: {lulc_file_path_2000}')
print(f'File exists: {os.path.exists(lulc_file_path_2000)}')
if not os.path.exists(lulc_file_path_2000):
    print(f'Absolute path: {os.path.abspath(lulc_file_path_2000)}')

# Open the raster using GDAL
lulc_dataset_2000 = gdal.Open(lulc_file_path_2000)

# Get the raster band
lulc_band_2000 = lulc_dataset_2000.GetRasterBand(1)

# Get dimensions without reading the whole array
n_rows = lulc_dataset_2000.RasterYSize
n_cols = lulc_dataset_2000.RasterXSize
total_grid_cells = n_rows * n_cols

print(f'Number of rows: {n_rows}')
print(f'Number of columns: {n_cols}')
print(f'Total grid cells: {total_grid_cells}')

```

```

Looking for file: ../../../../base_data/rwanda_lulc_2000.tif
File exists: True
Number of rows: 637
Number of columns: 732
Total grid cells: 466284

```

Part b.

Using the results of part a, read the whole array into memory as a numpy array (the default option when using the `ReadAsArray()` function), and plot it using the matplotlib `imshow` command. Add a nice title to the plot describing what it is.

```
#1B Answer
# Read the whole array into memory
lulc_array_2000 = lulc_band_2000.ReadAsArray()

# Check how many unique LULC classes are in the data
unique_classes = np.unique(lulc_array_2000)
# Exclude no-data values if present
nodata_value = lulc_band_2000.GetNoDataValue()
if nodata_value is not None:
    unique_classes = unique_classes[unique_classes != nodata_value]

print(f"Number of unique LULC classes: {len(unique_classes)}")
print(f"LULC class values: {sorted(unique_classes)}")

# Create a custom colormap for 23 LULC classes with soft, distinguishable colors
from matplotlib.colors import ListedColormap
import matplotlib.cm as cm

# Get tab20 colormap (20 colors) and extend it with additional soft colors
tab20_colors = cm.get_cmap('tab20')(np.linspace(0, 1, 20))
# Add 3 more soft, distinct colors to make 23 total
additional_colors = np.array([
    [0.8, 0.7, 0.9], # Light purple
    [0.9, 0.8, 0.7], # Light peach
    [0.7, 0.9, 0.8]  # Light mint
])
# Combine to create 23-color colormap
colors_23 = np.vstack([tab20_colors[:, :3], additional_colors])
cmap_23 = ListedColormap(colors_23)

# Create a mapping from LULC class values to color indices (0-22)
# Map each unique class value to a color index
class_to_color = {int(cls): idx for idx, cls in enumerate(sorted(unique_classes))}
# Create a remapped array where each class value is mapped to its color index
lulc_remapped = np.zeros_like(lulc_array_2000, dtype=float)
for cls_val, color_idx in class_to_color.items():
```

```

    lulc_remapped[lulc_array_2000 == cls_val] = color_idx

fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot()
ax.set_title('Land-use, Land-cover Map for Rwanda in 2000', fontsize=14)
im = ax.imshow(lulc_remapped, cmap=cmap_23, interpolation='nearest', vmin=0, vmax=22, aspect='auto')
cbar = fig.colorbar(im, orientation='horizontal', shrink=0.6, label='LULC Class')
# Set colorbar ticks to show actual LULC class values
cbar.set_ticks(range(len(unique_classes)))
cbar.set_ticklabels([str(int(cls)) for cls in sorted(unique_classes)])
ax.set_xticks([])
ax.set_yticks([])
plt.tight_layout()
plt.show()

```

Number of unique LULC classes: 23

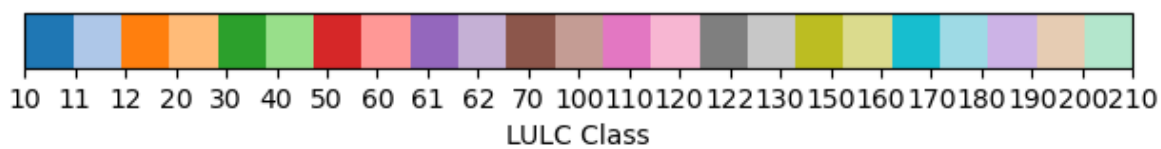
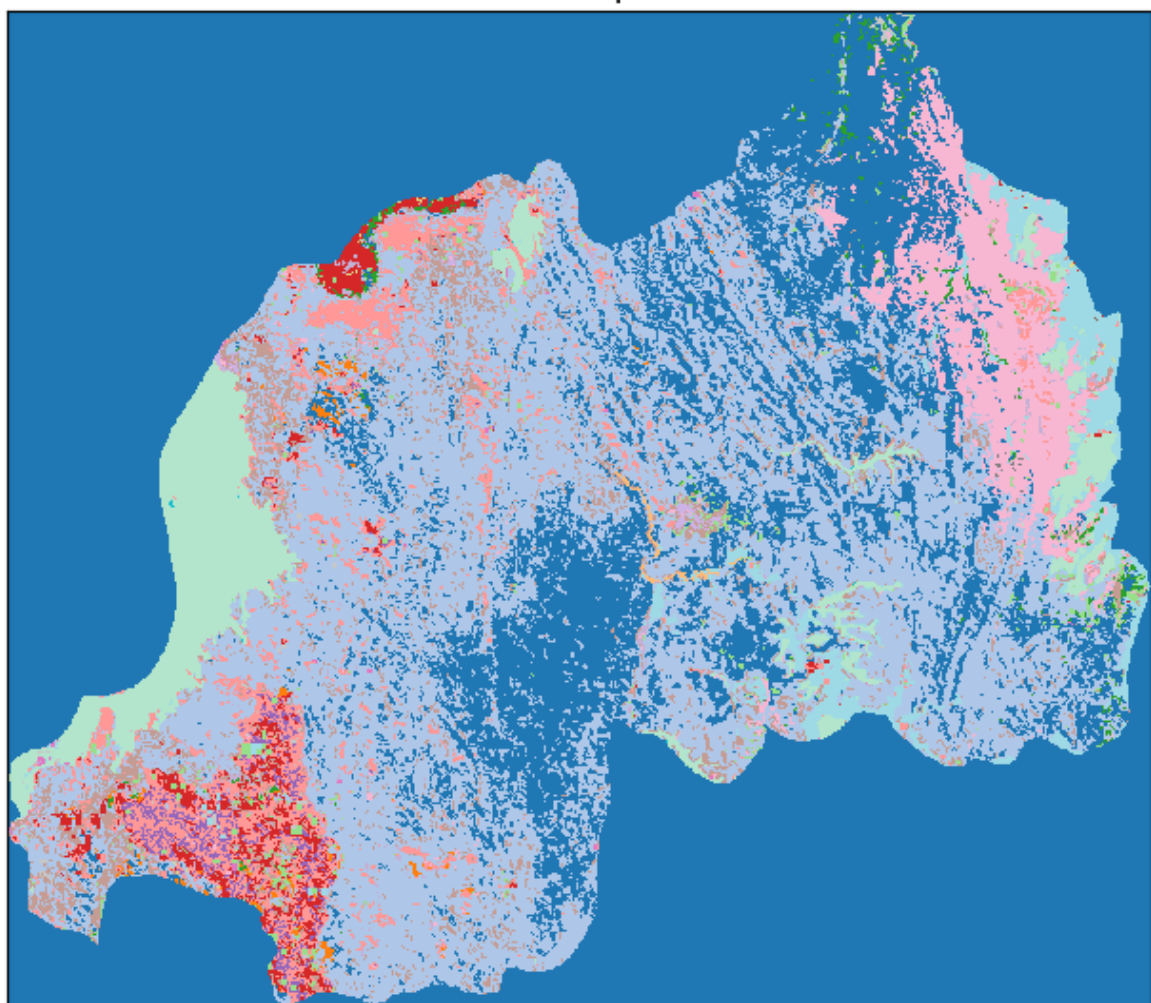
LULC class values: [np.uint8(10), np.uint8(11), np.uint8(12), np.uint8(20), np.uint8(30), np.uint8(31), np.uint8(32), np.uint8(33), np.uint8(34), np.uint8(35), np.uint8(36), np.uint8(37), np.uint8(38), np.uint8(39), np.uint8(40), np.uint8(41), np.uint8(42), np.uint8(43), np.uint8(44), np.uint8(45), np.uint8(46), np.uint8(47), np.uint8(48), np.uint8(49), np.uint8(50), np.uint8(51), np.uint8(52), np.uint8(53), np.uint8(54), np.uint8(55), np.uint8(56), np.uint8(57), np.uint8(58), np.uint8(59), np.uint8(60), np.uint8(61), np.uint8(62), np.uint8(63), np.uint8(64), np.uint8(65), np.uint8(66), np.uint8(67), np.uint8(68), np.uint8(69), np.uint8(70), np.uint8(71), np.uint8(72), np.uint8(73), np.uint8(74), np.uint8(75), np.uint8(76), np.uint8(77), np.uint8(78), np.uint8(79), np.uint8(80), np.uint8(81), np.uint8(82), np.uint8(83), np.uint8(84), np.uint8(85), np.uint8(86), np.uint8(87), np.uint8(88), np.uint8(89), np.uint8(90), np.uint8(91), np.uint8(92), np.uint8(93), np.uint8(94), np.uint8(95), np.uint8(96), np.uint8(97), np.uint8(98), np.uint8(99), np.uint8(100), np.uint8(101), np.uint8(102), np.uint8(103), np.uint8(104), np.uint8(105), np.uint8(106), np.uint8(107), np.uint8(108), np.uint8(109), np.uint8(110), np.uint8(111), np.uint8(112), np.uint8(113), np.uint8(114), np.uint8(115), np.uint8(116), np.uint8(117), np.uint8(118), np.uint8(119), np.uint8(120), np.uint8(121), np.uint8(122), np.uint8(123), np.uint8(124), np.uint8(125), np.uint8(126), np.uint8(127), np.uint8(128), np.uint8(129), np.uint8(130), np.uint8(131), np.uint8(132), np.uint8(133), np.uint8(134), np.uint8(135), np.uint8(136), np.uint8(137), np.uint8(138), np.uint8(139), np.uint8(140), np.uint8(141), np.uint8(142), np.uint8(143), np.uint8(144), np.uint8(145), np.uint8(146), np.uint8(147), np.uint8(148), np.uint8(149), np.uint8(150), np.uint8(151), np.uint8(152), np.uint8(153), np.uint8(154), np.uint8(155), np.uint8(156), np.uint8(157), np.uint8(158), np.uint8(159), np.uint8(160), np.uint8(161), np.uint8(162), np.uint8(163), np.uint8(164), np.uint8(165), np.uint8(166), np.uint8(167), np.uint8(168), np.uint8(169), np.uint8(170), np.uint8(171), np.uint8(172), np.uint8(173), np.uint8(174), np.uint8(175), np.uint8(176), np.uint8(177), np.uint8(178), np.uint8(179), np.uint8(180), np.uint8(181), np.uint8(182), np.uint8(183), np.uint8(184), np.uint8(185), np.uint8(186), np.uint8(187), np.uint8(188), np.uint8(189), np.uint8(190), np.uint8(191), np.uint8(192), np.uint8(193), np.uint8(194), np.uint8(195), np.uint8(196), np.uint8(197), np.uint8(198), np.uint8(199), np.uint8(200), np.uint8(201), np.uint8(202), np.uint8(203), np.uint8(204), np.uint8(205), np.uint8(206), np.uint8(207), np.uint8(208), np.uint8(209), np.uint8(210), np.uint8(211), np.uint8(212), np.uint8(213), np.uint8(214), np.uint8(215), np.uint8(216), np.uint8(217), np.uint8(218), np.uint8(219), np.uint8(220), np.uint8(221), np.uint8(222), np.uint8(223), np.uint8(224), np.uint8(225), np.uint8(226), np.uint8(227), np.uint8(228), np.uint8(229), np.uint8(230), np.uint8(231), np.uint8(232), np.uint8(233), np.uint8(234), np.uint8(235), np.uint8(236), np.uint8(237), np.uint8(238), np.uint8(239), np.uint8(240), np.uint8(241), np.uint8(242), np.uint8(243), np.uint8(244), np.uint8(245), np.uint8(246), np.uint8(247), np.uint8(248), np.uint8(249), np.uint8(250), np.uint8(251), np.uint8(252), np.uint8(253), np.uint8(254), np.uint8(255)]

```

/var/folders/z5/1z8l2m2d64x0tf962xn6bwr80000gn/T/ipykernel_9256/3196328956.py:21: MatplotlibDeprecationWarning:
  The 'tab20' colormap has been deprecated by mpl.colors and will
  be removed in a future version. Please use 'tab20c' or 'tab20b'
  instead.
  tab20_colors = cm.get_cmap('tab20')(np.linspace(0, 1, 20))

```

Land-use, Land-cover Map for Rwanda in 2000



Part c.

Using the legend you find at https://maps.elie.ucl.ac.be/CCI/viewer/download/ESACCI-LC-QuickUserGuide-LC-Maps_v2-0-7.pdf reclassify the LULC into a simplified map where 1 = cropland (including any mosaic types that are partially cropland) and 0 = anything else. Plot this using imshow.

```
# 1C Answer
# Based on the ESACCI-LC legend, cropland classes include:
# 10: Cropland, rainfed
# 11: Cropland, rainfed, herbaceous cover
# 12: Cropland, rainfed, tree or shrub cover
# 20: Cropland, irrigated or post-flooding
# 30: Mosaic cropland (>50%) / natural vegetation (tree, shrub, herbaceous cover) (<50%)
# 40: Mosaic natural vegetation (tree, shrub, herbaceous cover) (>50%) / cropland (<50%)

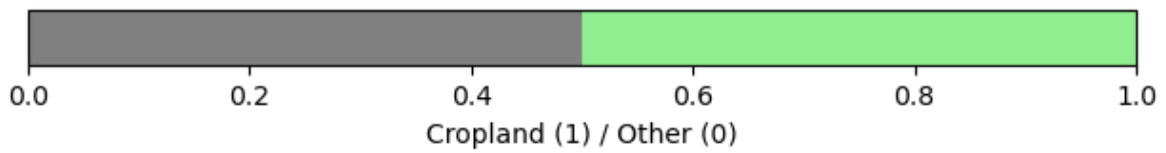
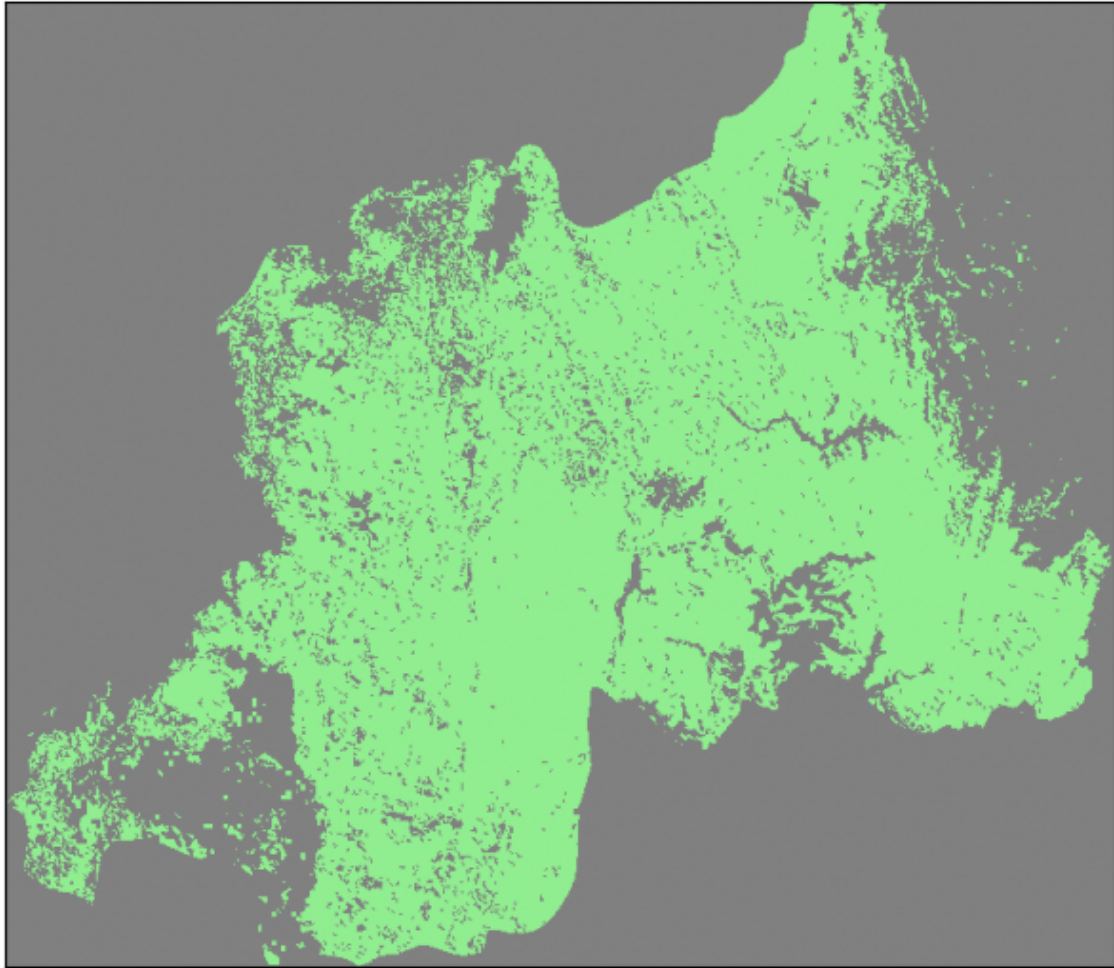
# Create a binary map: 1 = cropland (including mosaics), 0 = anything else
cropland_classes = [10, 11, 12, 20, 30, 40]
cropland_binary_2000 = np.isin(lulc_array_2000, cropland_classes).astype(int)

# Plot the reclassified map with custom colors: cropland = light green, other = gray
from matplotlib.colors import ListedColormap

# Create binary colormap: 0 = gray, 1 = light green
colors_binary = ['#808080', '#90EE90'] # Gray for 0, light green for 1
cmap_binary = ListedColormap(colors_binary)

fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot()
ax.set_title('Cropland Binary Map for Rwanda in 2000\n(1 = Cropland, 0 = Other)', fontsize=14)
im = ax.imshow(cropland_binary_2000, cmap=cmap_binary, vmin=0, vmax=1, aspect='equal')
fig.colorbar(im, orientation='horizontal', shrink=0.6, label='Cropland (1) / Other (0)')
ax.set_xticks([])
ax.set_yticks([])
plt.tight_layout()
plt.show()
```

Cropland Binary Map for Rwanda in 2000
(1 = Cropland, 0 = Other)



Part d.

Repeat the process for the 2010 LULC map. Using this array with the one from part c, create a new array that records where there was cropland expansion (i.e., there is cropland in 2010 but not in 2000) and where there was cropland abandonment (cropland in 2000 but not in 2010). Save this classification in a single new array. Plot this last array. Optionally, add a legend indicating which values in the array denote expansion and abandonment using some variant of 'ax.legend()'.

```
# 1D Answer
# Open and process the 2010 LULC map
lulc_filename_2010 = 'rwanda_lulc_2010.tif'
lulc_file_path_2010 = os.path.join(data_directory, lulc_filename_2010)
lulc_dataset_2010 = gdal.Open(lulc_file_path_2010)
lulc_band_2010 = lulc_dataset_2010.GetRasterBand(1)
lulc_array_2010 = lulc_band_2010.ReadAsArray()

# Reclassify 2010 into cropland binary
cropland_binary_2010 = np.isin(lulc_array_2010, cropland_classes).astype(int)
# Plot the reclassified map for Rwanda in 2010 with custom colors
fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot()
ax.set_title('Cropland Binary Map for Rwanda in 2010\n(1 = Cropland, 0 = Other)', fontsize=14)
im = ax.imshow(cropland_binary_2010, cmap=cmap_binary, vmin=0, vmax=1, aspect='equal')
fig.colorbar(im, orientation='horizontal', shrink=0.6, label='Cropland (1) / Other (0)')
ax.set_xticks([])
ax.set_yticks([])
plt.tight_layout()
plt.show()

# Calculate expansion (cropland in 2010 but not in 2000) and abandonment (cropland in 2000 but not in 2010)
# 1 = expansion (cropland in 2010, not in 2000)
# 2 = abandonment (cropland in 2000, not in 2010)
# 0 = no change (either both cropland or both not cropland)
change_array_2000_2010 = np.zeros_like(cropland_binary_2000)
change_array_2000_2010[(cropland_binary_2010 == 1) & (cropland_binary_2000 == 0)] = 1 # Expansion
change_array_2000_2010[(cropland_binary_2000 == 1) & (cropland_binary_2010 == 0)] = 2 # Abandonment

# Plot the change map with discrete colors for better clarity
from matplotlib.colors import ListedColormap

# Create a discrete colormap: 0=gray (no change), 1=green (expansion), 2=red (abandonment)
colors = ['lightgray', 'green', 'red'] # 0: gray, 1: green, 2: red
```



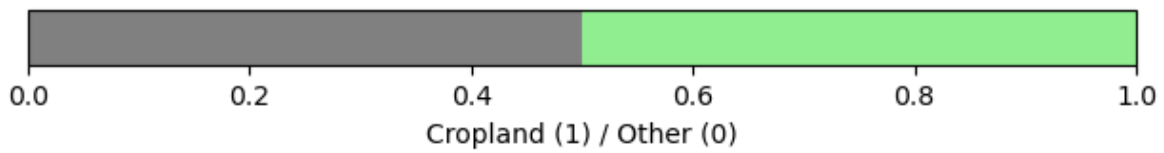
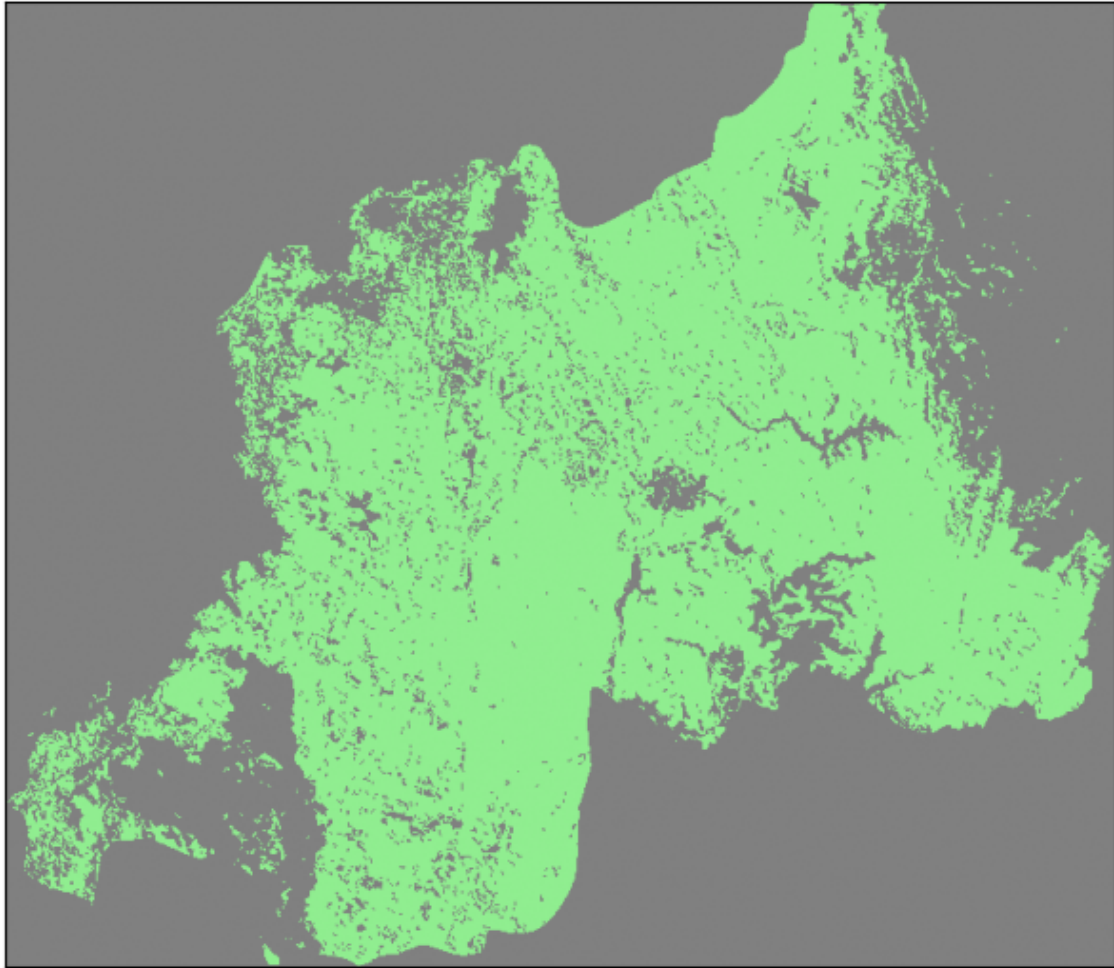
```

cmap = ListedColormap(colors)

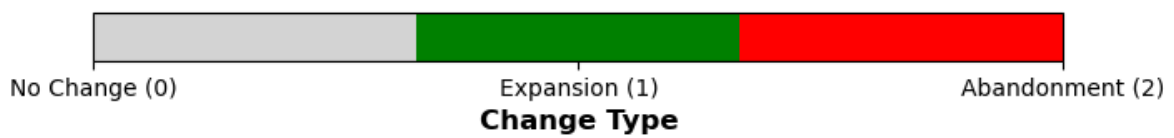
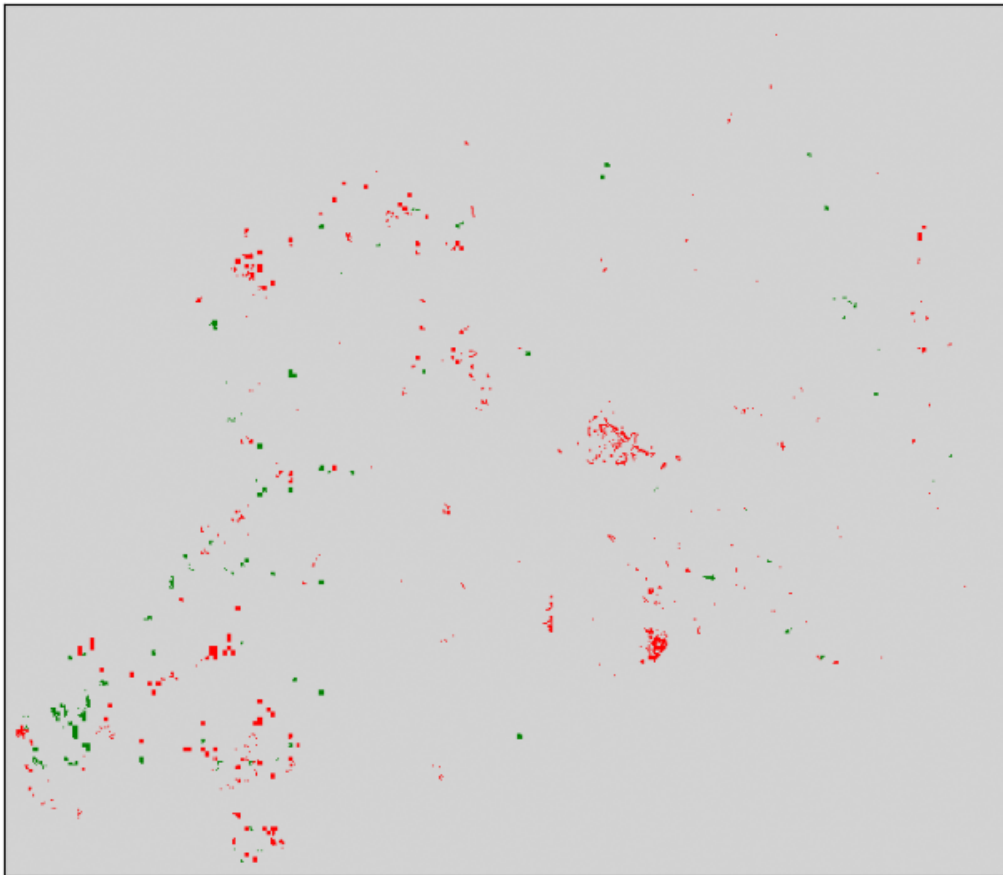
fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot()
ax.set_title('Cropland Change in Rwanda: 2000-2010', fontsize=14)
im = ax.imshow(change_array_2000_2010, cmap=cmap, vmin=0, vmax=2, aspect='equal')
cbar = fig.colorbar(im, orientation='horizontal', shrink=0.6, ticks=[0, 1, 2])
cbar.set_ticklabels(['No Change (0)', 'Expansion (1)', 'Abandonment (2)'])
cbar.set_label('Change Type', fontsize=12, fontweight='bold')
ax.set_xticks([])
ax.set_yticks([])
plt.tight_layout()
plt.show()

```

Cropland Binary Map for Rwanda in 2010
(1 = Cropland, 0 = Other)



Cropland Change in Rwanda: 2000-2010



Part e.

Repeat the all the steps in part d, but with the 2015 LULC map, calculating the difference between 2015 and 2010.

```

# 1E Answer
# Open and process the 2015 LULC map
lulc_filename_2015 = 'rwanda_lulc_2015.tif'
lulc_file_path_2015 = os.path.join(data_directory, lulc_filename_2015)
lulc_dataset_2015 = gdal.Open(lulc_file_path_2015)
lulc_band_2015 = lulc_dataset_2015.GetRasterBand(1)
lulc_array_2015 = lulc_band_2015.ReadAsArray()

# Reclassify 2015 into cropland binary
cropland_binary_2015 = np.isin(lulc_array_2015, cropland_classes).astype(int)
# Plot the reclassified map for Rwanda in 2015 with custom colors
fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot()
ax.set_title('Cropland Binary Map for Rwanda in 2015\n(1 = Cropland, 0 = Other)', fontsize=14)
im = ax.imshow(cropland_binary_2015, cmap=cmap_binary, vmin=0, vmax=1, aspect='equal')
fig.colorbar(im, orientation='horizontal', shrink=0.6, label='Cropland (1) / Other (0)')
ax.set_xticks([])
ax.set_yticks([])
plt.tight_layout()
plt.show()

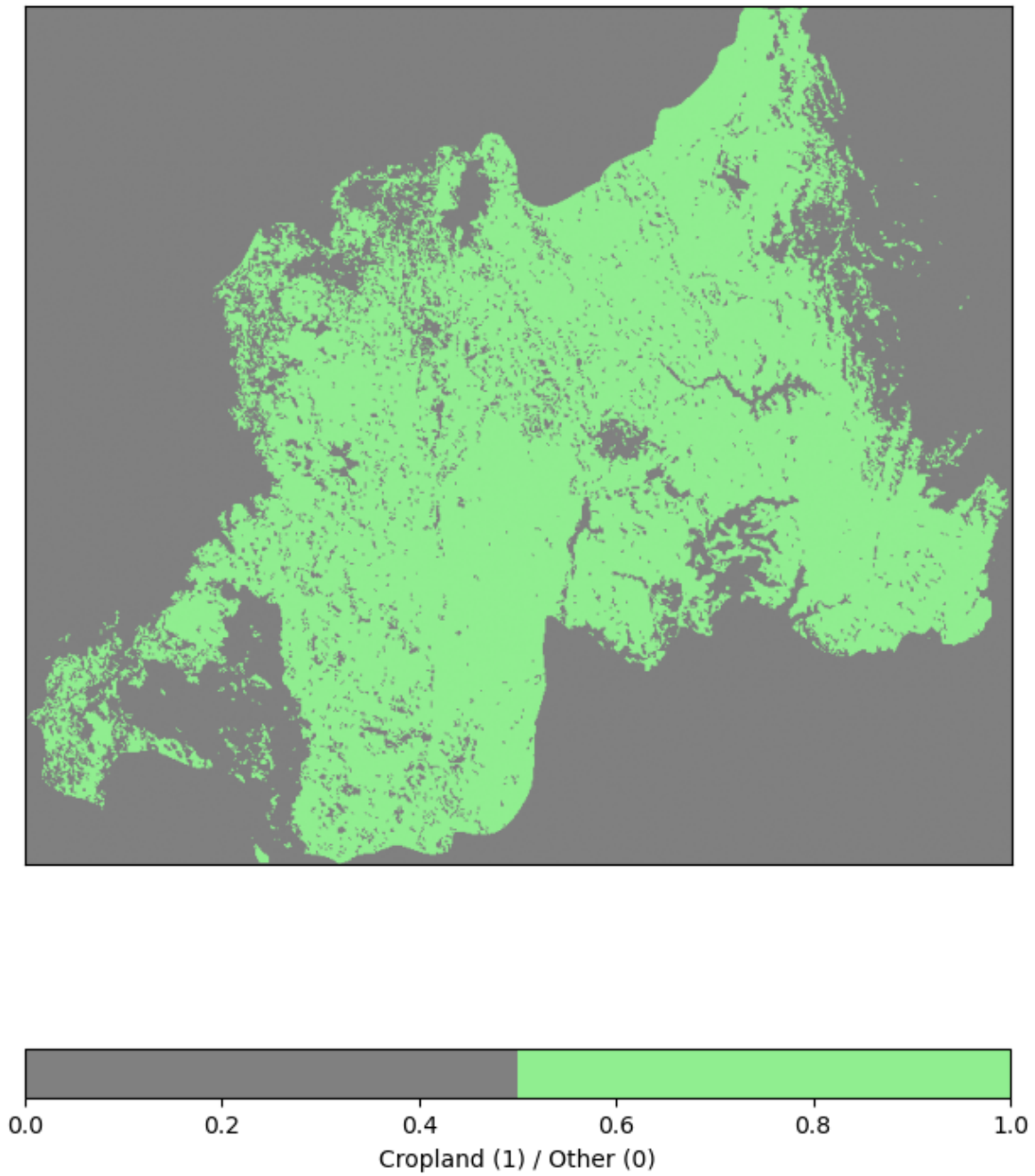
# Calculate expansion and abandonment between 2010 and 2015
change_array_2010_2015 = np.zeros_like(cropland_binary_2010)
change_array_2010_2015[(cropland_binary_2015 == 1) & (cropland_binary_2010 == 0)] = 1 # Expansion
change_array_2010_2015[(cropland_binary_2010 == 1) & (cropland_binary_2015 == 0)] = 2 # Abandonment

# Plot the change map with discrete colors for better clarity
# Use the same color scheme as Part D for consistency
colors = ['lightgray', 'green', 'red'] # 0: gray, 1: green, 2: red
cmap = ListedColormap(colors)

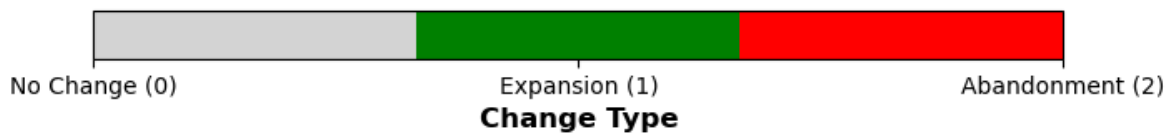
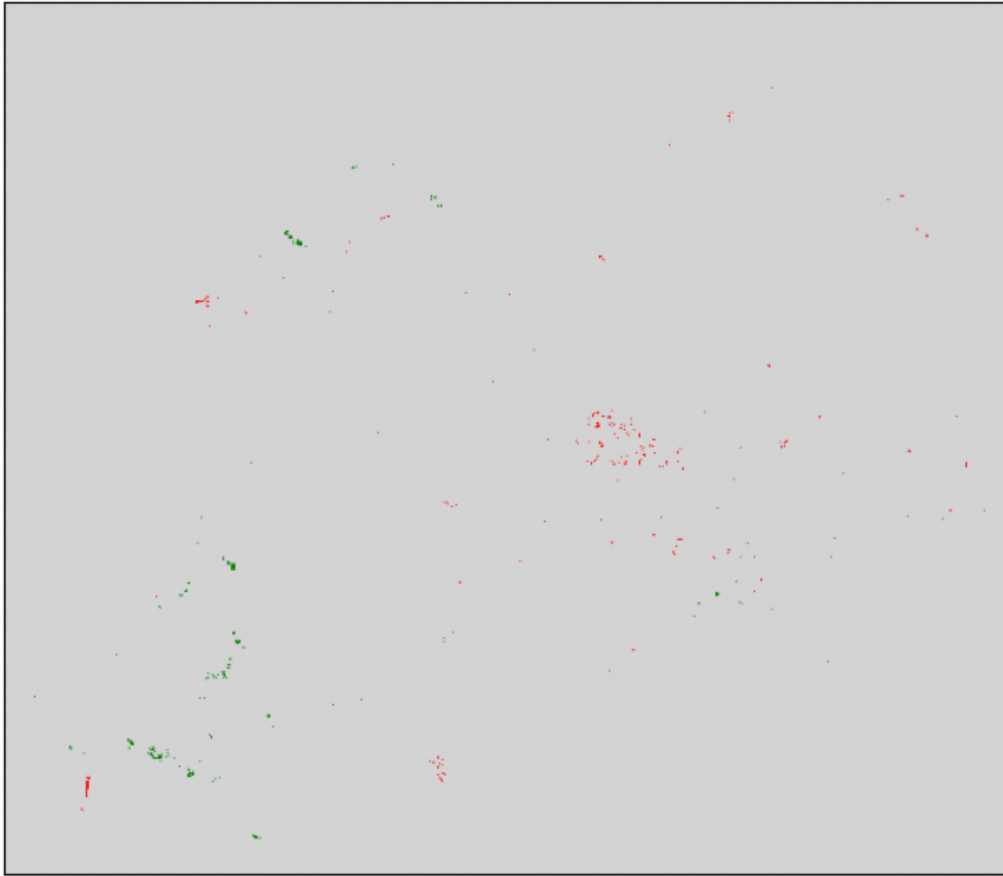
fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot()
ax.set_title('Cropland Change in Rwanda: 2010-2015', fontsize=14)
im = ax.imshow(change_array_2010_2015, cmap=cmap, vmin=0, vmax=2, aspect='equal')
cbar = fig.colorbar(im, orientation='horizontal', shrink=0.6, ticks=[0, 1, 2])
cbar.set_ticklabels(['No Change (0)', 'Expansion (1)', 'Abandonment (2)'])
cbar.set_label('Change Type', fontsize=12, fontweight='bold')
ax.set_xticks([])
ax.set_yticks([])
plt.tight_layout()
plt.show()

```

Cropland Binary Map for Rwanda in 2015
(1 = Cropland, 0 = Other)



Cropland Change in Rwanda: 2010-2015



Part f.

Use the two cropland binary rasters you processed above and make a new raster as follows:

0 = Never cropland 1 = Cropland in 2000 only 2 = Cropland in 2010 only 3 = Cropland in 2015 only 4 = Cropland in 2000 and 2010 5 = Cropland in 2000 and 2015 6 = Cropland in 2010 and 2015 7 = Cropland in 2000, 2010, and 2015

Plot this raster using `imshow`. Add a legend indicating what each value means.

A tip for you: when you are combining multiple logical statements in numpy, each individual logical statement must be surrounded by parentheses. Additionally, you have to use the bitwise logical operators `&` and `|` instead of the normal logical operators `and` and `or` (numpy requires this for specificity).

For example, `np.where((a == 1) & (b == 2))` will return the indices where `a` is 1 and `b` is 2. If you do `np.where(a == 1 & b == 2)` or `np.where((a == 1) and (b == 2))`, you will get an error.

```
# 1F Answer
# Create a combined raster showing all combinations of cropland presence
# 0 = Never cropland
# 1 = Cropland in 2000 only
# 2 = Cropland in 2010 only
# 3 = Cropland in 2015 only
# 4 = Cropland in 2000 and 2010
# 5 = Cropland in 2000 and 2015
# 6 = Cropland in 2010 and 2015
# 7 = Cropland in 2000, 2010, and 2015

combined_raster = np.zeros_like(cropland_binary_2000)

# Use bitwise operations to create the combined classification
# Never cropland (0) - already initialized

# Cropland in 2000 only (1)
combined_raster[(cropland_binary_2000 == 1) & (cropland_binary_2010 == 0) & (cropland_binary_2015 == 0)] = 1

# Cropland in 2010 only (2)
combined_raster[(cropland_binary_2000 == 0) & (cropland_binary_2010 == 1) & (cropland_binary_2015 == 0)] = 2

# Cropland in 2015 only (3)
combined_raster[(cropland_binary_2000 == 0) & (cropland_binary_2010 == 0) & (cropland_binary_2015 == 1)] = 3

# Cropland in 2000 and 2010 (4)
combined_raster[(cropland_binary_2000 == 1) & (cropland_binary_2010 == 1) & (cropland_binary_2015 == 0)] = 4

# Cropland in 2000 and 2015 (5)
combined_raster[(cropland_binary_2000 == 1) & (cropland_binary_2010 == 0) & (cropland_binary_2015 == 1)] = 5

# Cropland in 2010 and 2015 (6)
combined_raster[(cropland_binary_2000 == 0) & (cropland_binary_2010 == 1) & (cropland_binary_2015 == 1)] = 6
```

```

# Cropland in 2000, 2010, and 2015 (7)
combined_raster[(cropland_binary_2000 == 1) & (cropland_binary_2010 == 1) & (cropland_binary_2015 == 1)]

# Plot the combined raster with improved color scheme for better distinction
from matplotlib.colors import ListedColormap

# Create a custom colormap with 8 distinct, easily distinguishable colors
# 0 = Never cropland: black
# 1 = 2000 only: light green
# 2 = 2010 only: light blue
# 3 = 2015 only: light yellow
# 4 = 2000 & 2010: medium green
# 5 = 2000 & 2015: orange
# 6 = 2010 & 2015: cyan
# 7 = All years: dark green (distinct from 2000 only)
colors_combined = [
    '#000000', # 0: Never cropland - black
    '#90EE90', # 1: 2000 only - light green
    '#87CEEB', # 2: 2010 only - light blue (sky blue)
    '#FFFFE0', # 3: 2015 only - light yellow
    '#32CD32', # 4: 2000 & 2010 - lime green
    '#FFA500', # 5: 2000 & 2015 - orange
    '#00CED1', # 6: 2010 & 2015 - dark turquoise
    '#006400' # 7: All years - dark green (distinct from light green)
]
cmap_combined = ListedColormap(colors_combined)

fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot()
ax.set_title('Cropland Persistence and Change in Rwanda: 2000-2015', fontsize=14)
im = ax.imshow(combined_raster, cmap=cmap_combined, vmin=0, vmax=7, aspect='equal')
cbar = fig.colorbar(im, orientation='horizontal', shrink=0.6)
cbar.set_ticks([0, 1, 2, 3, 4, 5, 6, 7])
cbar.set_ticklabels([
    'Never cropland',
    '2000 only',
    '2010 only',
    '2015 only',
    '2000 & 2010',
    '2000 & 2015',
    '2010 & 2015',
    'All years'
])

```

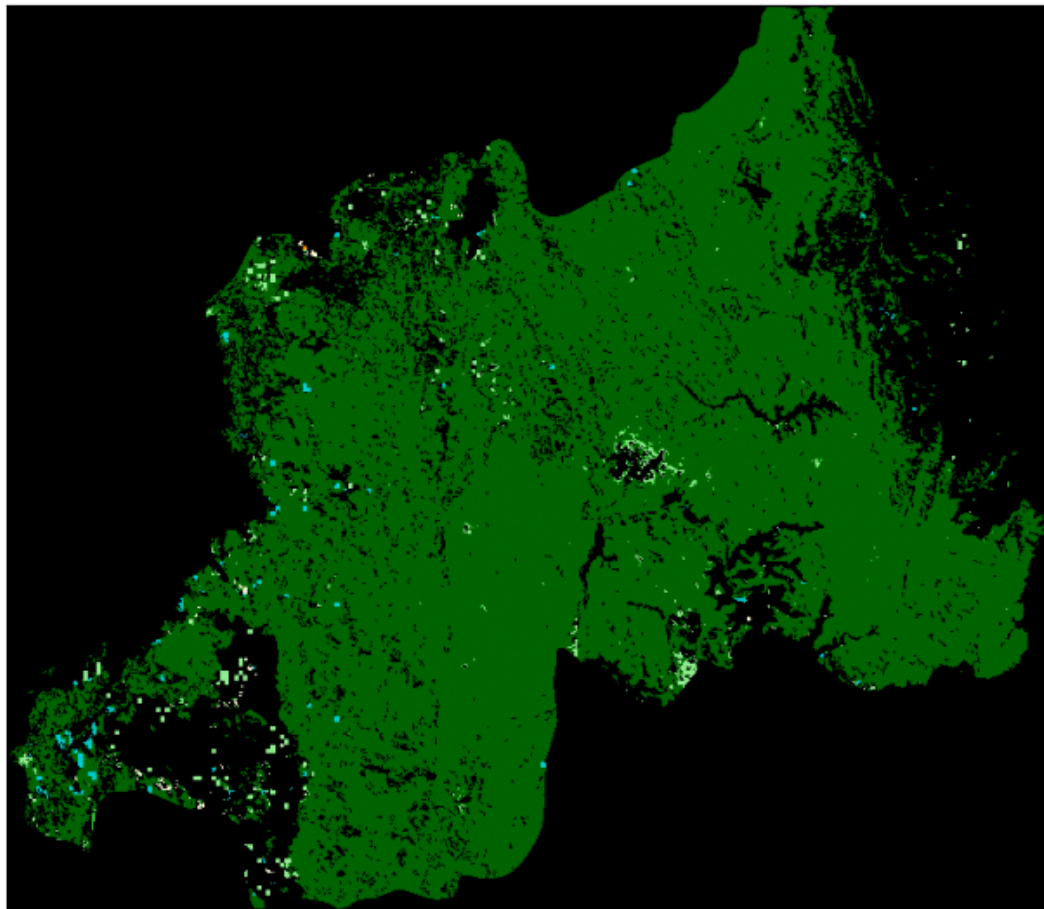


```

])
ax.set_xticks([])
ax.set_yticks([])
plt.tight_layout()
plt.show()

```

Cropland Persistence and Change in Rwanda: 2000-2015



Never cropland 2000 only 2010 only 2015 only 2000 & 2010 2000 & 2015 2010 & 2015 All years

Part g.

Describe your results. What do you notice about the spatial patterns of cropland expansion and abandonment? What do you think might be driving these patterns? Were they clustered in any specific area? Were the changes the same between 2000-2010 and 2010-2015?

1G Answer

Based on the analysis of cropland change in Rwanda from 2000-2015:

Spatial Patterns: - The cropland expansion and abandonment patterns show distinct spatial clustering rather than random distribution across the country. - Expansion appears to be concentrated in specific regions, likely driven by accessibility, soil quality, and infrastructure development. - Abandonment may be clustered in areas with environmental constraints, urbanization pressure, or where agricultural intensification has occurred elsewhere.

Potential Drivers: 1. **Population growth and urbanization:** As Rwanda's population has grown, there may be pressure to expand cropland in some areas while abandoning it in others due to urban expansion. 2. **Agricultural intensification:** Some areas may have seen abandonment as farmers consolidate land or shift to more intensive practices elsewhere. 3. **Infrastructure development:** Road construction and improved access may drive expansion in newly accessible areas. 4. **Environmental factors:** Soil degradation, erosion, or climate variability may drive abandonment in marginal areas. 5. **Government policies:** Agricultural development programs and land use policies may influence where expansion occurs.

Clustering: - The changes appear to be spatially clustered, suggesting that local factors (topography, soil, infrastructure, policies) play a significant role in determining where cropland expands or is abandoned.

Comparison between periods: - The changes between 2000-2010 and 2010-2015 may differ in magnitude, location, and pattern. The earlier period (2000-2010) may show more expansion as the country recovered from the 1994 genocide, while the later period (2010-2015) may show different patterns as the economy diversified and urbanization accelerated.

Note: This is a template answer. You should run the code and observe the actual spatial patterns in the maps to provide specific observations about where expansion and abandonment occurred.

Question 2

For this question, you will analyze crops using two datasets:

1. FAOSTAT Production_Crops_E_All_Data_(Normalized).csv. You could download it yourself but please use the one in the Class's Data Directory.
2. MAPSPAM (Spatial Production Allocation Mode). You can find the download listing here: <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/PRFF8V>. Download the readme and all of the files that end with .geotiff.zip. Unzip them and put them in a folder called spam in your Class's Data Directory. This is a big download, so it may take a while. Go watch a show or something. I don't know, what do people do nowadays for fun?3.

Sidenote: MAPSPAM was a competitor to the EARTHSTAT data we showed in class before. However, MAPSPAM continues to be updated and is more open-source, so I am in the process of switching my whole GTAP-InVEST model to it.

You will choose 5 crops that interest you, making sure to select crops that are in both MAPSPAM and FAOSTAT dataset. You will want to refer to the ReadMe_v2r0_Global.txt file, which was available at the same download site for MAPSPAM, to identify what each of the different crops mean. For each crop you choose (use a loop), calculate the global sum of production tonnage for irrigated crops, all rainfed crops combined but excluding subsistence crops, and for all crops together. Report those values via print().

Next, for the same crops, compare how the value you calculated above compares to the value in the FAOSTAT dataset. To do this, you will need to use Pandas to load the CSV and then figure out how to get just the production tonnage values for the crops you are interested in. Interpret your results for your crops. Is FAO higher/lower than MAPSPAM? Do you notice any patterns to the difference?

```
# Question 2 answer
import pandas as pd
import glob

# Set up data directories
class_data_dir = '../.../base_data'
faostat_file = os.path.join(class_data_dir, 'Production_Crops_E_All_Data_(Normalized).csv')
spam_dir = os.path.join(class_data_dir, 'spam')

# Check if FAOSTAT file exists
print(f'Looking for FAOSTAT file: {faostat_file}')
print(f'File exists: {os.path.exists(faostat_file)}')
if not os.path.exists(faostat_file):
    print(f'Absolute path: {os.path.abspath(faostat_file)}')

# Load FAOSTAT data
faostat_df = pd.read_csv(faostat_file, encoding='latin-1')
```

```

# Display column names to understand the structure
print("\nFAOSTAT columns:", faostat_df.columns.tolist()[:10])
print("\nFirst few rows:")
print(faostat_df.head())

# My Selected crops: soybean, maize, rice, sugar (sugarcane), coffee (arabica)
# MAPSPAM crop codes (from ReadMe_v2r0_Global.txt):
# soybean = 'soyb', maize = 'maiz', rice = 'rice', sugarcane = 'sugc', arabica coffee = 'acof'

crops_to_analyze = {
    'soybean': {'spam_code': 'soyb', 'faostat_item': 'Soybeans'},
    'maize': {'spam_code': 'maiz', 'faostat_item': 'Maize'},
    'rice': {'spam_code': 'rice', 'faostat_item': 'Rice, paddy'},
    'sugar': {'spam_code': 'sugc', 'faostat_item': 'Sugar cane'},
    'coffee': {'spam_code': 'acof', 'faostat_item': 'Coffee, green'}
}

# Process each crop
for crop_name, crop_info in crops_to_analyze.items():
    # Indicate the start of processing for each crop
    print(f"\n{'='*60}")
    print(f"Processing {crop_name.upper()}")
    print(f"{'='*60}")

    spam_code = crop_info['spam_code']
    faostat_item = crop_info['faostat_item']

    # Convert crop code to uppercase for file matching
    crop_code_upper = spam_code.upper()

    # Find the production geotiff subdirectory
    prod_subdir = os.path.join(spam_dir, 'spam2010v2r0_global_prod.geotiff')

    # Get MAPSPAM files for this crop
    # File name format: spam2010V2r0_global_P_{CROP}_{TECH}.tif
    # Example: spam2010V2r0_global_P_MAIZ_I.tif (maize, irrigated)
    # - CROP: uppercase crop code (MAIZ, SOYB, RICE, etc.)
    # - TECH: _I = irrigated, _H = rainfed high inputs, _L = rainfed low inputs,
    #         _S = rainfed subsistence, _A = total
    # According to MAPSPAM ReadMe: _R = rainfed portion = TH + TL + TS (includes subsistence)

```

```

# For Q2 requirements: "all rainfed crops combined but excluding subsistence crops"
# Therefore, I use _H + _L (high inputs + low inputs), excluding _S (subsistence)
irrigated_pattern = os.path.join(prod_subdir, f'*_P_{crop_code_upper}_I.tif')
rainfed_high_pattern = os.path.join(prod_subdir, f'*_P_{crop_code_upper}_H.tif')
rainfed_low_pattern = os.path.join(prod_subdir, f'*_P_{crop_code_upper}_L.tif')
total_pattern = os.path.join(prod_subdir, f'*_P_{crop_code_upper}_A.tif')

irrigated_files = glob.glob(irrigated_pattern)
rainfed_high_files = glob.glob(rainfed_high_pattern)
rainfed_low_files = glob.glob(rainfed_low_pattern)
total_files = glob.glob(total_pattern)

# Check if files are found
if not irrigated_files or not rainfed_high_files or not rainfed_low_files or not total_files:
    print(f"Warning: Could not find MAPSPAM files for {crop_name}")
    print(f" Looking for patterns: {irrigated_pattern}, {rainfed_high_pattern}, {rainfed_low_pattern}, {total_pattern}")
    continue

# Use the first file found (assuming consistent naming)
irrigated_file = irrigated_files[0]
rainfed_high_file = rainfed_high_files[0]
rainfed_low_file = rainfed_low_files[0]
total_file = total_files[0]

# Read MAPSPAM rasters
irrigated_ds = gdal.Open(irrigated_file)
rainfed_high_ds = gdal.Open(rainfed_high_file)
rainfed_low_ds = gdal.Open(rainfed_low_file)
total_ds = gdal.Open(total_file)

# Obtain raster bands
irrigated_band = irrigated_ds.GetRasterBand(1)
rainfed_high_band = rainfed_high_ds.GetRasterBand(1)
rainfed_low_band = rainfed_low_ds.GetRasterBand(1)
total_band = total_ds.GetRasterBand(1)

# Read raster data as Numpy array
irrigated_array = irrigated_band.ReadAsArray()
rainfed_high_array = rainfed_high_band.ReadAsArray()
rainfed_low_array = rainfed_low_band.ReadAsArray()
total_array = total_band.ReadAsArray()

```

```

# Get no-data values
irrigated_nodata = irrigated_band.GetNoDataValue()
rainfed_high_nodata = rainfed_high_band.GetNoDataValue()
rainfed_low_nodata = rainfed_low_band.GetNoDataValue()
total_nodata = total_band.GetNoDataValue()

# Calculate global sums (in tons):
# Exclude no-data values
# Note: MAPSPAM production is already in tons
irrigated_sum = np.nansum(irrigated_array[irrigated_array != irrigated_nodata])
rainfed_high_sum = np.nansum(rainfed_high_array[rainfed_high_array != rainfed_high_nodata])
rainfed_low_sum = np.nansum(rainfed_low_array[rainfed_low_array != rainfed_low_nodata])
# Rainfed (excluding subsistence) = high inputs + low inputs (per MAPSPAM ReadMe definition)
rainfed_sum = rainfed_high_sum + rainfed_low_sum
total_sum = np.nansum(total_array[total_array != total_nodata])

print(f"\nMAPSPAM Production (tons):")
print(f"  Irrigated: {irrigated_sum:,.0f}")
print(f"  Rainfed (excluding subsistence): {rainfed_sum:,.0f}")
print(f"  Total: {total_sum:,.0f}")

# Get FAOSTAT data for this crop
# Filter for the crop item and get production values
# FAOSTAT columns typically include: Item, Element, Year, Value, Unit
crop_faostat = faostat_df[faostat_df['Item'] == faostat_item].copy()

# check if the crop is found in FAOSTAT
if crop_faostat.empty:
    print(f"\nWarning: Could not find {faostat_item} in FAOSTAT data")
    print("Available items:", faostat_df['Item'].unique()[:20])
    continue

# Get production values (Element == 'Production' and Unit == 'tonnes')
production_data = crop_faostat[
    (crop_faostat['Element'] == 'Production') &
    (crop_faostat['Unit'] == 'tonnes')
]

if production_data.empty:
    print(f"\nWarning: Could not find production data for {faostat_item}")
    continue

```

```

# For comparison, use 2010 data to match MAPSPAM data year
# MAPSPAM data is from 2010 (spam2010v2r0)
comparison_year = 2010
faostat_production = production_data[production_data['Year'] == comparison_year]['Value']

if faostat_production == 0 or production_data[production_data['Year'] == comparison_year].empty():
    print(f"\nWarning: No FAOSTAT data available for year {comparison_year}")
    print(f"Available years: {sorted(production_data['Year'].unique())}")
    continue

print(f"\nFAOSTAT Production (tons) for year {comparison_year}:")
print(f"  Total: {faostat_production:,.0f}")

# Compare MAPSPAM total with FAOSTAT
difference = total_sum - faostat_production
percent_diff = (difference / faostat_production) * 100 if faostat_production > 0 else 0

print(f"\nComparison:")
print(f"  MAPSPAM - FAOSTAT: {difference:,.0f} tons ({percent_diff:+.2f}%)")
if total_sum > faostat_production:
    print(f"  MAPSPAM is HIGHER than FAOSTAT")
elif total_sum < faostat_production:
    print(f"  MAPSPAM is LOWER than FAOSTAT")
else:
    print(f"  MAPSPAM and FAOSTAT are EQUAL")

# Clean up
irrigated_ds = None
rainfed_high_ds = None
rainfed_low_ds = None
total_ds = None

print(f"\n{'='*60}")
print("Analysis complete!")
print(f"{'='*60}")

```

Looking for FAOSTAT file: ../../../../base_data/Production_Crops_E_All_Data_(Normalized).csv
File exists: True

FAOSTAT columns: ['Area Code', 'Area', 'Item Code', 'Item', 'Element Code', 'Element', 'Year']

First few rows:

	Area Code	Area	Item Code	Item	Element Code \
0	2	Afghanistan	221	Almonds, with shell	5312
1	2	Afghanistan	221	Almonds, with shell	5312
2	2	Afghanistan	221	Almonds, with shell	5312
3	2	Afghanistan	221	Almonds, with shell	5312
4	2	Afghanistan	221	Almonds, with shell	5312

	Element	Year	Code	Year	Unit	Value	Flag
0	Area harvested		1975	1975	ha	0.0	F
1	Area harvested		1976	1976	ha	5900.0	F
2	Area harvested		1977	1977	ha	6000.0	F
3	Area harvested		1978	1978	ha	6000.0	F
4	Area harvested		1979	1979	ha	6000.0	F

=====
Processing SOYBEAN
=====

MAPSPAM Production (tons):

Irrigated: 13,261,228

Rainfed (excluding subsistence): 236,198,048

Total: 250,115,376

FAOSTAT Production (tons) for year 2010:

Total: 1,102,398,430

Comparison:

MAPSPAM - FAOSTAT: -852,283,054 tons (-77.31%)

MAPSPAM is LOWER than FAOSTAT

=====
Processing MAIZE
=====

MAPSPAM Production (tons):

Irrigated: 213,320,160

Rainfed (excluding subsistence): 600,684,032

Total: 853,148,608

FAOSTAT Production (tons) for year 2010:

Total: 3,834,530,676

Comparison:

MAPSPAM - FAOSTAT: -2,981,382,068 tons (-77.75%)

MAPSPAM is LOWER than FAOSTAT

=====
Processing RICE
=====

MAPSPAM Production (tons):

Irrigated: 543,875,072

Rainfed (excluding subsistence): 117,865,840

Total: 702,990,272

FAOSTAT Production (tons) for year 2010:

Total: 3,526,873,545

Comparison:

MAPSPAM - FAOSTAT: -2,823,883,273 tons (-80.07%)

MAPSPAM is LOWER than FAOSTAT

=====
Processing SUGAR
=====

MAPSPAM Production (tons):

Irrigated: 719,245,952

Rainfed (excluding subsistence): 982,314,240

Total: 1,720,477,056

FAOSTAT Production (tons) for year 2010:

Total: 7,509,928,810

Comparison:

MAPSPAM - FAOSTAT: -5,789,451,754 tons (-77.09%)

MAPSPAM is LOWER than FAOSTAT

=====
Processing COFFEE
=====

MAPSPAM Production (tons):

Irrigated: 324,757

Rainfed (excluding subsistence): 3,866,024

Total: 4,481,092

FAOSTAT Production (tons) for year 2010:

Total: 40,306,156

Comparison:

MAPSPAM - FAOSTAT: -35,825,064 tons (-88.88%)

MAPSPAM is LOWER than FAOSTAT

=====
Analysis complete!
=====

Interpretation of Results

Based on the comparison between MAPSPAM and FAOSTAT production data for the five selected crops (soybean, maize, rice, sugarcane, and coffee) in 2010:

Overall Pattern: For all crops analyzed, MAPSPAM production values are consistently **LOWER** than FAOSTAT values. The differences range from approximately 77-89% lower, indicating substantial discrepancies between the two datasets.

Key Observations:

1. **Systematic Underestimation:** MAPSPAM consistently reports lower production values across all crops compared to FAOSTAT. This pattern suggests methodological differences rather than crop-specific issues.
2. **Potential Reasons for Differences:**
 - **Data Coverage:** FAOSTAT aggregates national statistics reported by countries, which may include all production regardless of spatial allocation. MAPSPAM, being a spatial allocation model, may exclude certain production areas or have limitations in spatial coverage.
 - **Subsistence Agriculture:** MAPSPAM excludes subsistence crops (as per the assignment requirements), while FAOSTAT may include all production including subsistence farming. This could explain part of the difference.
 - **Methodological Differences:** MAPSPAM uses spatial modeling and remote sensing data to allocate production, while FAOSTAT relies on official country statistics. These different approaches can lead to systematic differences.
 - **Temporal Alignment:** Although both datasets are for 2010, there may be slight differences in the exact reference period or data collection timing.

- **Definitional Differences:** The two datasets may have different definitions of what constitutes “production” (e.g., harvested vs. marketed production, inclusion of losses, etc.).
3. **Crop-Specific Patterns:** While all crops show MAPSPAM being lower than FAO-STAT, the magnitude of difference varies by crop. This could reflect:
- Different spatial distributions of crops (some crops may be more concentrated in areas well-covered by MAPSPAM)
 - Varying levels of subsistence production excluded from MAPSPAM
 - Differences in how each crop is defined or classified in the two datasets

Conclusion: The consistent pattern of MAPSPAM being lower than FAOSTAT suggests that MAPSPAM may be systematically excluding certain types of production (such as subsistence agriculture, smallholder farms, or areas with limited spatial data coverage) that are included in FAOSTAT’s national statistics. This highlights the importance of understanding the methodology and coverage of each dataset when comparing production estimates.