

Jinal Vyas 1233053785 Part - 1 Data Preparation Report

```

1 # Here, I am using the data of bands 3 to 7 from Landsat-8 of NASA over the Alaska Mendenhall Glacier
2
3 # The product I am using is Landsat 8 OLI/TIRS Collection 2 atmospherically corrected surface reflectance
4 # the raw reflectance values can be influenced by the time of day, weather conditions, and the atmospheric conditions
5
6 # For the EDA representation purpose, I am taking the data of two months from 1st October 2023 to 1st November 2023
7 # Landsat collects images every 16 days, but here I have derived the images after applying cloud mask
8 # I have performed EDA on the tif file I derived from the median of pixels from two months of the above data
9 # In the next stage I will replicate this processing on data files of more years and perform segmentation

```

```

1 # Rasterio is a Python library used for reading, writing, and processing raster data, typically in GeoTIFF format
2 !pip install rasterio

```

```

Collecting rasterio
  Downloading rasterio-1.4.2-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (9.1 kB)
Collecting affine (from rasterio)
  Downloading affine-2.4.0-py3-none-any.whl.metadata (4.0 kB)
Requirement already satisfied: attrs in /usr/local/lib/python3.10/dist-packages (from rasterio) (24.2.0)
Requirement already satisfied: certifi in /usr/local/lib/python3.10/dist-packages (from rasterio) (2024.7.4)
Requirement already satisfied: click>=4.0 in /usr/local/lib/python3.10/dist-packages (from rasterio) (8.1.7)
Collecting cligj>=0.5 (from rasterio)
  Downloading cligj-0.7.2-py3-none-any.whl.metadata (5.0 kB)
Requirement already satisfied: numpy>=1.24 in /usr/local/lib/python3.10/dist-packages (from rasterio) (1.26.4)
Collecting click-plugins (from rasterio)
  Downloading click_plugins-1.1.1-py2.py3-none-any.whl.metadata (6.4 kB)
Requirement already satisfied: pyparsing in /usr/local/lib/python3.10/dist-packages (from rasterio) (3.2.1)
Downloading rasterio-1.4.2-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (22.2 MB)
  22.2/22.2 MB 23.1 MB/s eta 0:00:00
Downloading cligj-0.7.2-py3-none-any.whl (7.1 kB)
Downloading affine-2.4.0-py3-none-any.whl (15 kB)
Downloading click_plugins-1.1.1-py2.py3-none-any.whl (7.5 kB)
Installing collected packages: cligj, click-plugins, affine, rasterio
Successfully installed affine-2.4.0 click-plugins-1.1.1 cligj-0.7.2 rasterio-1.4.2

```

```

1 # geemap is the Python API of Google Earth Engine
2 import ee
3 import geemap.core as geemap
4 import torch
5 import rasterio
6 from torch.utils.data import Dataset, DataLoader
7 from torchvision import transforms
8 import numpy as np
9 import matplotlib.pyplot as plt
10 import numpy as np
11 from osgeo import gdal
12 import seaborn as sns

```

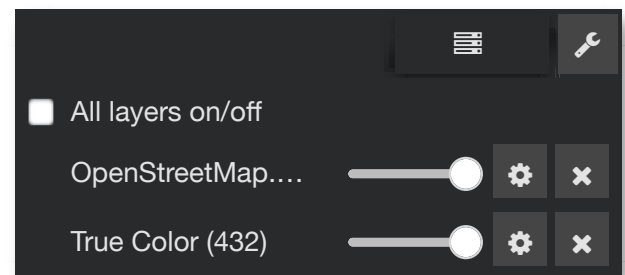
```

1 ee.Authenticate()
2 ee.Initialize(project='ee-jinalvyasict19')

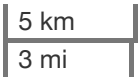
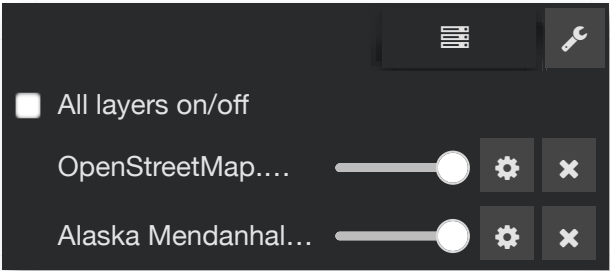
```

The image is a screenshot of a web application interface, likely a map viewer. On the left side, there is a vertical toolbar with four icons: a plus sign (+), a minus sign (-), a square with a dashed border, and a black square. Below these icons is a vertical stack of six black squares. At the bottom left, there is a scale bar with two segments, labeled "1000 km" and "500 mi". On the right side, there is a dark-themed layer control panel. It has a title bar with a list icon and a wrench icon. Below the title bar, there is a checkbox labeled "All layers on/off". Underneath, there are two layer entries: "OpenStreetMap..." and "True Color (432)". Each entry has a slider control, a gear icon for settings, and an 'x' icon for removal. The main area of the interface is a map of the United States, which is mostly obscured by the toolbar and the layer control panel. The map shows the outlines of the states in a light gray color.

```
1 # # viewing data over entire global map keeping Alaska Mendanhall Glacier as center
2 m = geemap.Map()
3 m.set_center(-134.62, 58.42, 6)
4 m.add_layer(dataset, visualization, 'True Color (432)')
5 m
```



```
1 # defining Mendanhall Glacier as region of interest
2 alaska_bounds = ee.Geometry.Polygon([[-134.62, 58.42], [-134.62, 58.59], [-134.42, 58.59], [-134.42,
3 map = geemap.Map()
4 map.centerObject(alaska_bounds, 10)
5 map.addLayer(alaska_bounds, {'color': 'lightblue'}, 'Alaska Mendanhall Glacier')
6 map
```



ipyleaflet | © [OpenStreetMap](#) contributors, Google Earth Engine

```
1 # If your dataset is made up of multiple images, you're creating a composite where the pixel value :
2 # from all the images at that location.
3 image = dataset.median().clip(alaska_bounds)
```

```
1 # downloading the image, am using high resolution of 10, which will increase data size a bit.
2 export_task = ee.batch.Export.image.toDrive(
3     image = image,
4     description = 'alaska_image',
5     scale = 10,
6     region = alaska_bounds,
7     maxPixels = 1e13,
8     fileFormat = 'GeoTIFF',
9     folder = 'EarthEngineExports_Alaska'
10 )
11
12 export_task.start()
```

```
1 # I will be using SR_B3 to SR_B7 for now, as they are enough for snow segmentation, other bands repi
2 # like a fusion of panchromatic and lower-resolution bands for better spatial resolution
3 image
```

► Image (19 bands)

[illegible]

- ✧ Exploratory Data Analysis

Mounted at /content/drive

```
{ 'AREA_OR_POINT': 'Area' }
Number of Bands: 19
GeoTransform: (-134.62000134216848, 8.983152841195215e-05, 0.0, 58.59000929237184, 0.0, -8.98315284119521
```

```
1 # making sure that the data is not entirely empty, hence loaded properly.
2 dataset.GetRasterBand(1).ReadAsArray()
```

```
array([[33328., 33328., 33328., ..., 44065., 44065., 43669.],
       [33328., 33328., 33328., ..., 44547., 44547., 44164.],
       [36721., 36721., 36721., ..., 44547., 44547., 44164.],
       ...,
       [ 7658.,  7658.,  7658., ..., 16850., 16850., 16850.],
       [ 7568.,  7568.,  7568., ..., 23831., 23831., 23831.],
       [ 7568.,  7568.,  7568., ..., 23831., 23831., 23831.]])
```

```
1 np.nanmax(dataset.GetRasterBand(1).ReadAsArray())
```

```
65454.0
```

```
1 # Printing the Meta-Data In Detail
2 with rasterio.open('/content/drive/My Drive/EarthEngineExports_Alaska/alaska_image.tif') as src:
3     print(f"Number of Bands: {src.count}")
4     for band in range(1, src.count + 1):
5         print(f"Band {band}: {src.descriptions[band - 1]}")
```

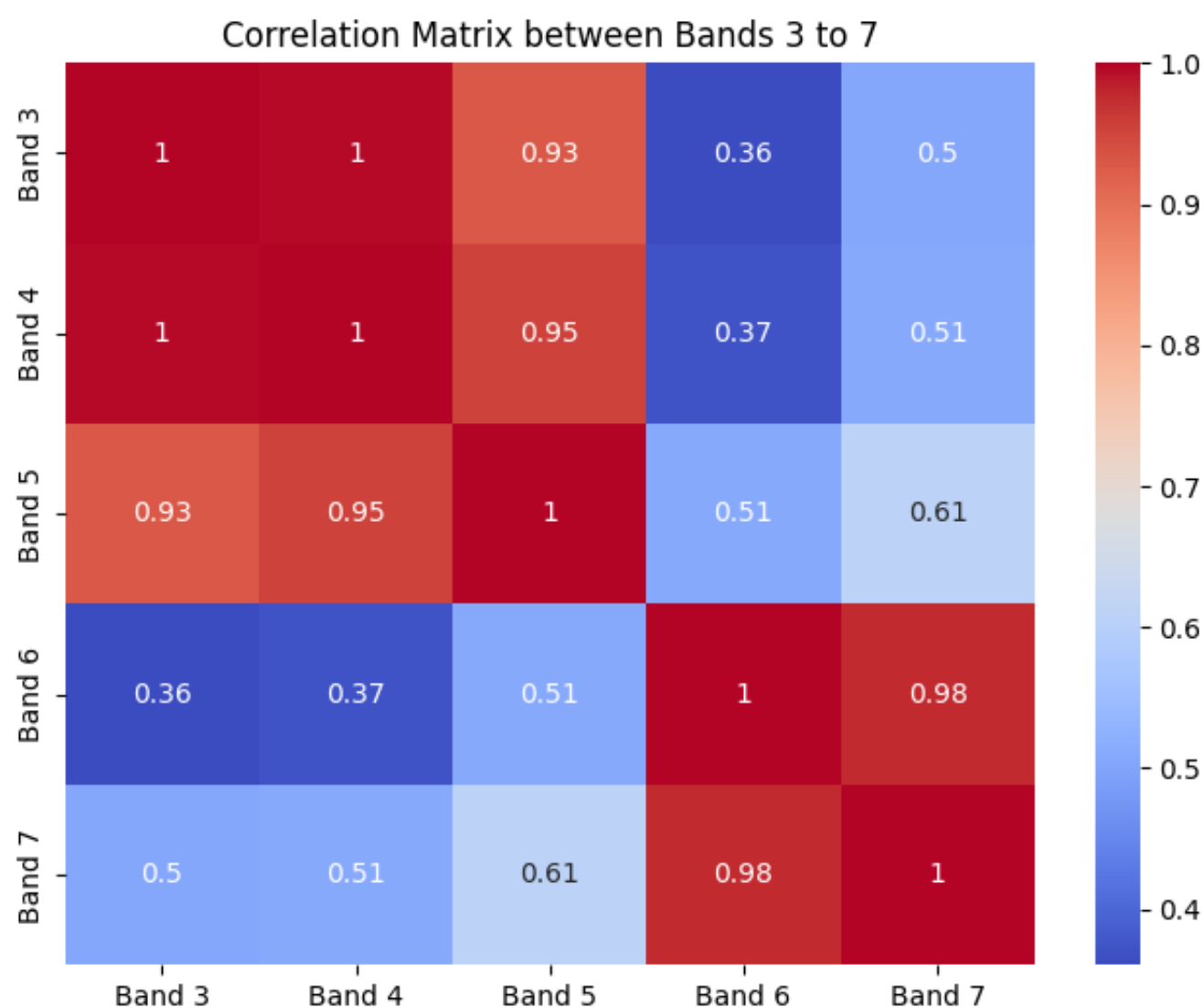
```
Number of Bands: 19
Band 1: SR_B1
Band 2: SR_B2
Band 3: SR_B3
Band 4: SR_B4
Band 5: SR_B5
Band 6: SR_B6
Band 7: SR_B7
Band 8: SR_QA_AEROSOL
Band 9: ST_B10
Band 10: ST_ATRAN
Band 11: ST_CDIST
Band 12: ST_DRAD
Band 13: ST_EMIS
Band 14: ST_EMSD
Band 15: ST_QA
Band 16: ST_TRAD
Band 17: ST_URAD
Band 18: QA_PIXEL
Band 19: QA_RADSAT
```

✓ Data Cleaning and EDA Continued

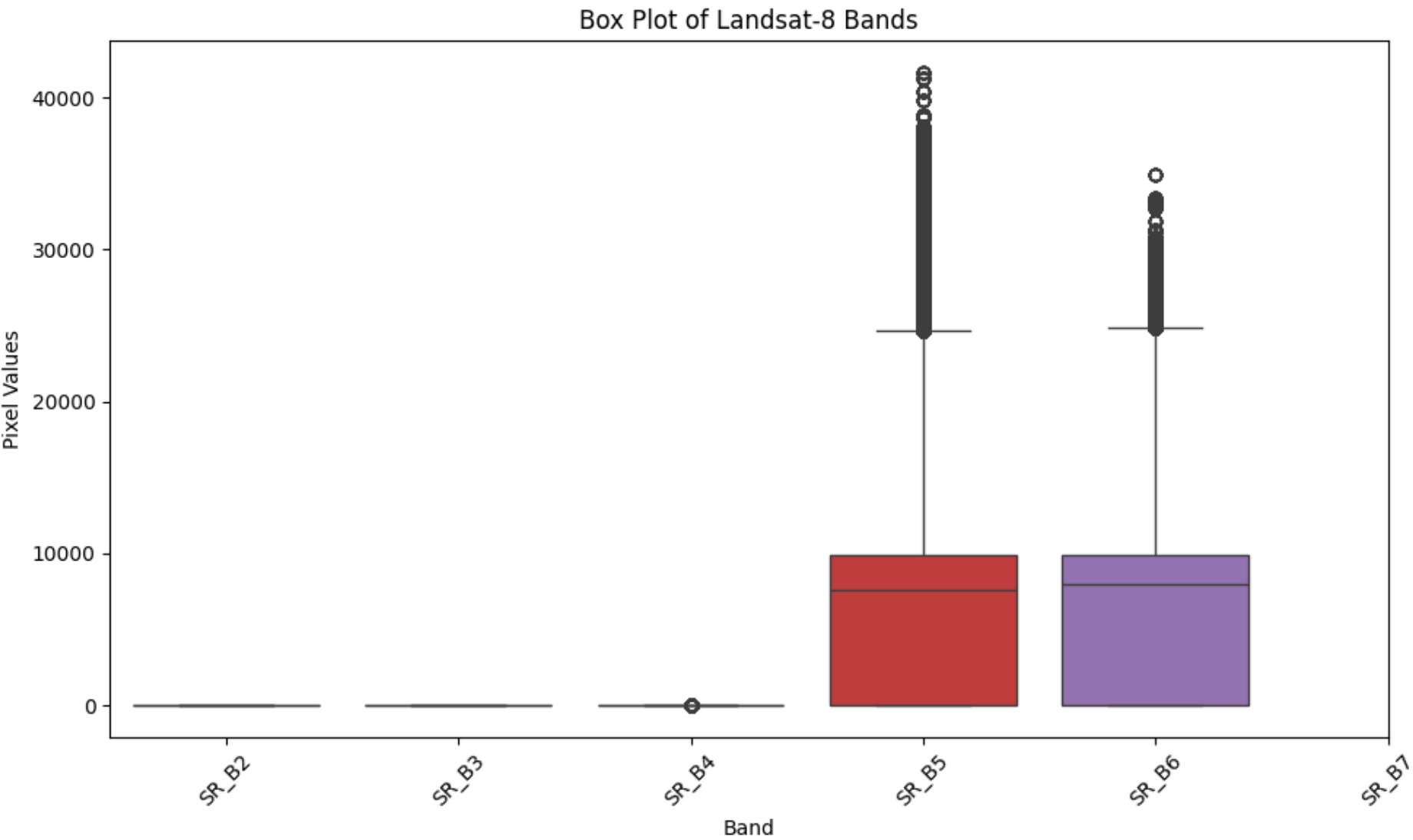
```

1 # Here, there are several Nan values because of these reasons:
2 # Satellite imagery can have NaN values in regions affected by clouds, cloud shadows, or other atmo:
3 # cannot be accurately measured.
4 # Some pixels may represent areas where no data is available (e.g., out of the image bounds, in rem:
5 # Hence, it is in best interest to replace these NaN values by just zero for simplicity of computat:
6 # values with arbitrary values (e.g., zero or the mean of surrounding pixels) introduces bias. It d:
7 # any further analysis, such as classification, index calculation (like NDSI or NDVI), or any machi:
8
9 # Plotting heatmap correlation matrix, to see relation between multiple bands, some pairs like 6-7,
10
11 with rasterio.open(file_path) as src:
12     band_3 = src.read(3)
13     band_4 = src.read(4)
14     band_5 = src.read(5)
15     band_6 = src.read(6)
16     band_7 = src.read(7)
17
18 band_3 = np.nan_to_num(band_3, nan=0)
19 band_4 = np.nan_to_num(band_4, nan=0)
20 band_5 = np.nan_to_num(band_5, nan=0)
21 band_6 = np.nan_to_num(band_6, nan=0)
22 band_7 = np.nan_to_num(band_7, nan=0)
23
24 band_3_flat = band_3.flatten()
25 band_4_flat = band_4.flatten()
26 band_5_flat = band_5.flatten()
27 band_6_flat = band_6.flatten()
28 band_7_flat = band_7.flatten()
29
30
31 data = np.stack([band_3_flat, band_4_flat, band_5_flat, band_6_flat, band_7_flat])
32
33 correlation_matrix = np.corrcoef(data)
34
35 # Create a heatmap of the correlation matrix
36 plt.figure(figsize=(8, 6))
37 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', xticklabels=["Band 3", "Band 4", "Band
38     yticklabels=["Band 3", "Band 4", "Band 5", "Band 6", "Band 7"])
39 plt.title('Correlation Matrix between Bands 3 to 7')
40 plt.show()

```



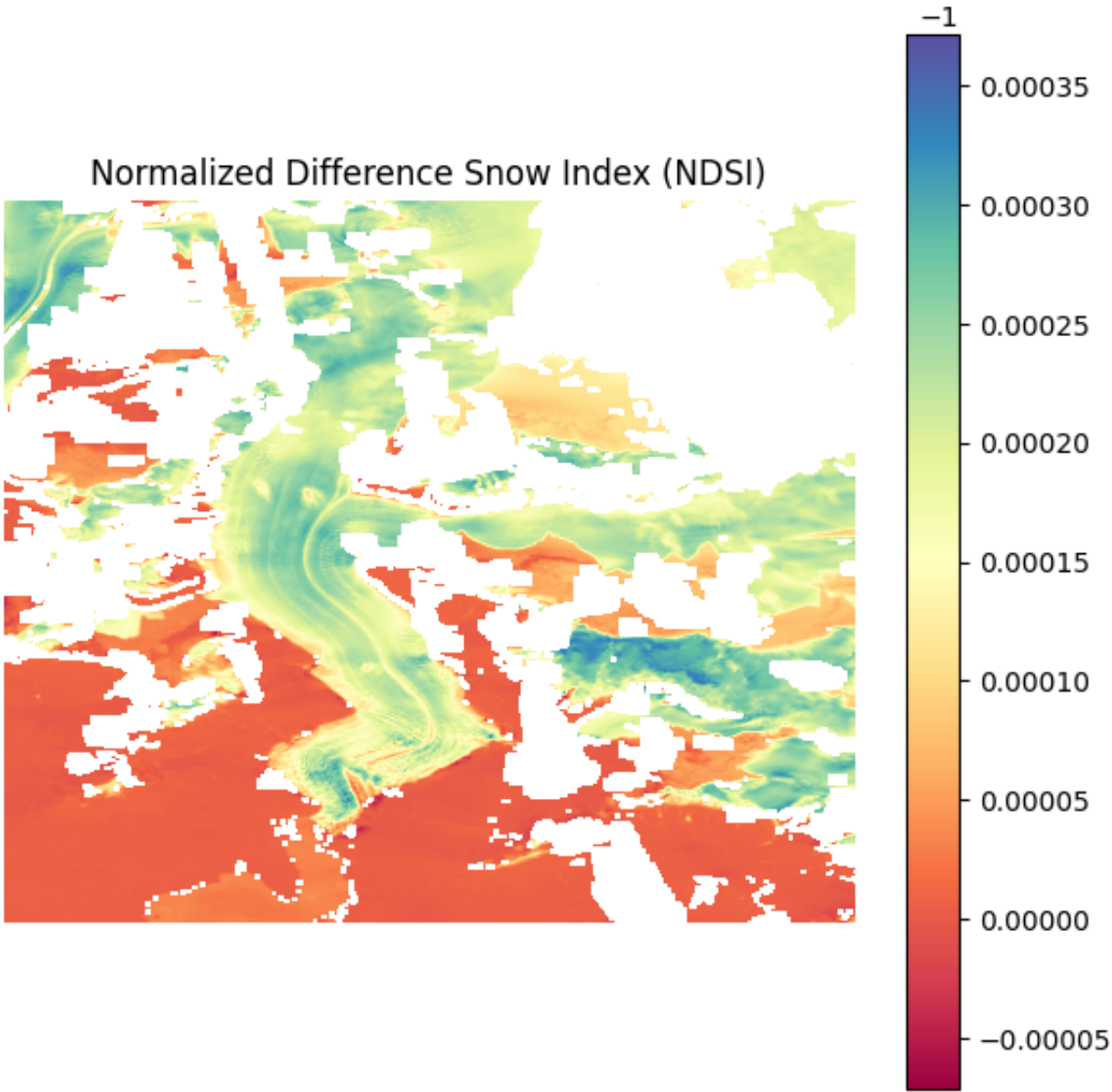

```
1 # Here, though there are outliers, but they may represent extreme temperature components like snow,
2 all_bands = [band_3_flat, band_4_flat, band_5_flat, band_6_flat, band_7_flat]
3 band_names = ['SR_B2', 'SR_B3', 'SR_B4', 'SR_B5', 'SR_B6', 'SR_B7']
4
5 plt.figure(figsize=(10, 6))
6 sns.boxplot(data=all_bands)
7 plt.xticks(np.arange(len(band_names)), band_names, rotation=45)
8 plt.title("Box Plot of Landsat-8 Bands")
9 plt.xlabel("Band")
10 plt.ylabel("Pixel Values")
11 plt.tight_layout()
12 plt.show()
```



Feature Engineering

```
1 # Performing Feature Engineering, i.e. deriving a new feature named NDSI (Normalized Difference Snow
2 # NDSI is typically used to detect snow and ice, as these features often have a high reflectance in
3 # a lower reflectance in the SWIR bands, which yields positive NDSI values
4
5 NDSI = (band_3 - band_6) / (band_3 + band_6)
6
7 # Plot NDSI
8 plt.figure(figsize=(7, 7))
9 plt.imshow(NDSI, cmap='Spectral')
10 plt.title('Normalized Difference Snow Index (NDSI)')
11 plt.colorbar()
12 plt.axis('off')
13 plt.show()
```

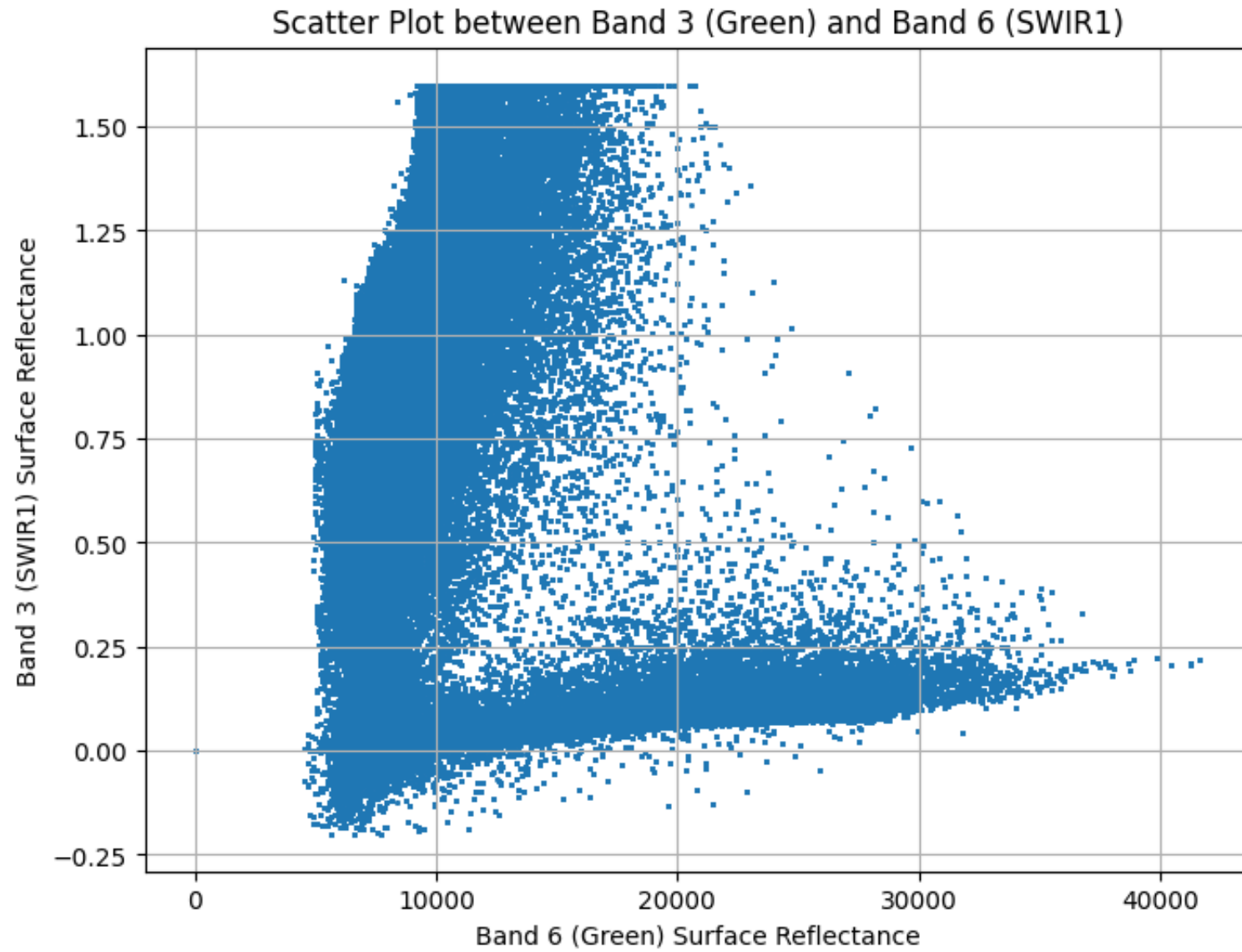
<ipython-input-185-277dda5d0dc9>:3: RuntimeWarning: invalid value encountered in divide
NDSI = (band_3 - band_6) / (band_3 + band_6)



```

1 # visualizing correlation between 3 and 6 to understand NDSI further
2 # For Snow/Ice: The Green band (Band 3) will have higher values (as snow reflects visible light), and
3 # (since snow absorbs infrared light). This results in a positive NDSI values.
4 plt.figure(figsize=(8, 6))
5 plt.scatter(band_6_flat, band_3_flat, alpha=0.5, s=1) # s=1 for small points, alpha for transparency
6 plt.title("Scatter Plot between Band 3 (Green) and Band 6 (SWIR1)")
7 plt.xlabel("Band 6 (Green) Surface Reflectance")
8 plt.ylabel("Band 3 (SWIR1) Surface Reflectance")
9 plt.grid(True)
10 plt.show()

```

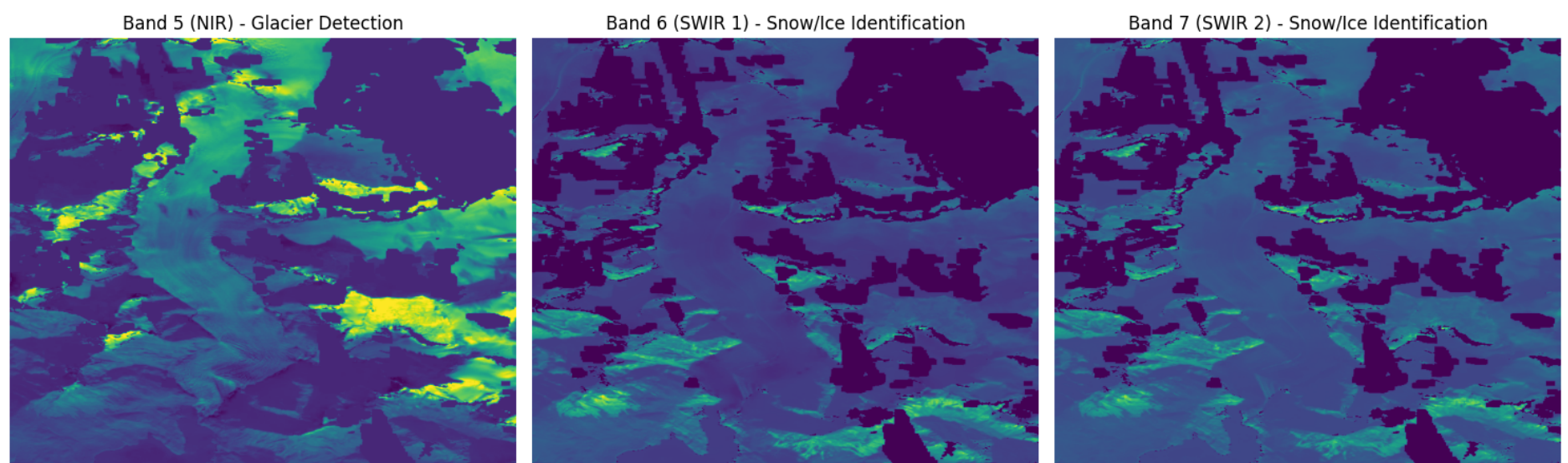


✓ EDA Continued

```

1 # Plotting band 5, 6 and 7.
2 # NIR (Band-5) is useful for vegetation monitoring, SWIR1 (Band-6) is sensitive to moisture content
3 # water tend to reflect differently in this band. SWIR2 (Band-7) is important for detecting moisture
4 # distinguishing water from other land covers.
5 file_path = "/content/drive/My Drive/EarthEngineExports_Alaska/alaska_image.tif"
6
7 with rasterio.open(file_path) as src:
8     band_5 = src.read(5) # NIR (Near Infrared)
9     band_6 = src.read(6) # SWIR 1
10    band_7 = src.read(7) # SWIR 2
11
12 band_5 = np.nan_to_num(band_5, nan=0)
13 band_6 = np.nan_to_num(band_6, nan=0)
14 band_7 = np.nan_to_num(band_7, nan=0)
15
16 fig, axes = plt.subplots(1, 3, figsize=(15, 5))
17
18 axes[0].imshow(band_5)
19 axes[0].set_title('Band 5 (NIR) - Glacier Detection')
20 axes[0].axis('off')
21
22 axes[1].imshow(band_6)
23 axes[1].set_title('Band 6 (SWIR 1) - Snow/Ice Identification')
24 axes[1].axis('off')
25
26
27 axes[2].imshow(band_7)
28 axes[2].set_title('Band 7 (SWIR 2) - Snow/Ice Identification')
29 axes[2].axis('off')
30
31 plt.tight_layout()
32 plt.show()

```

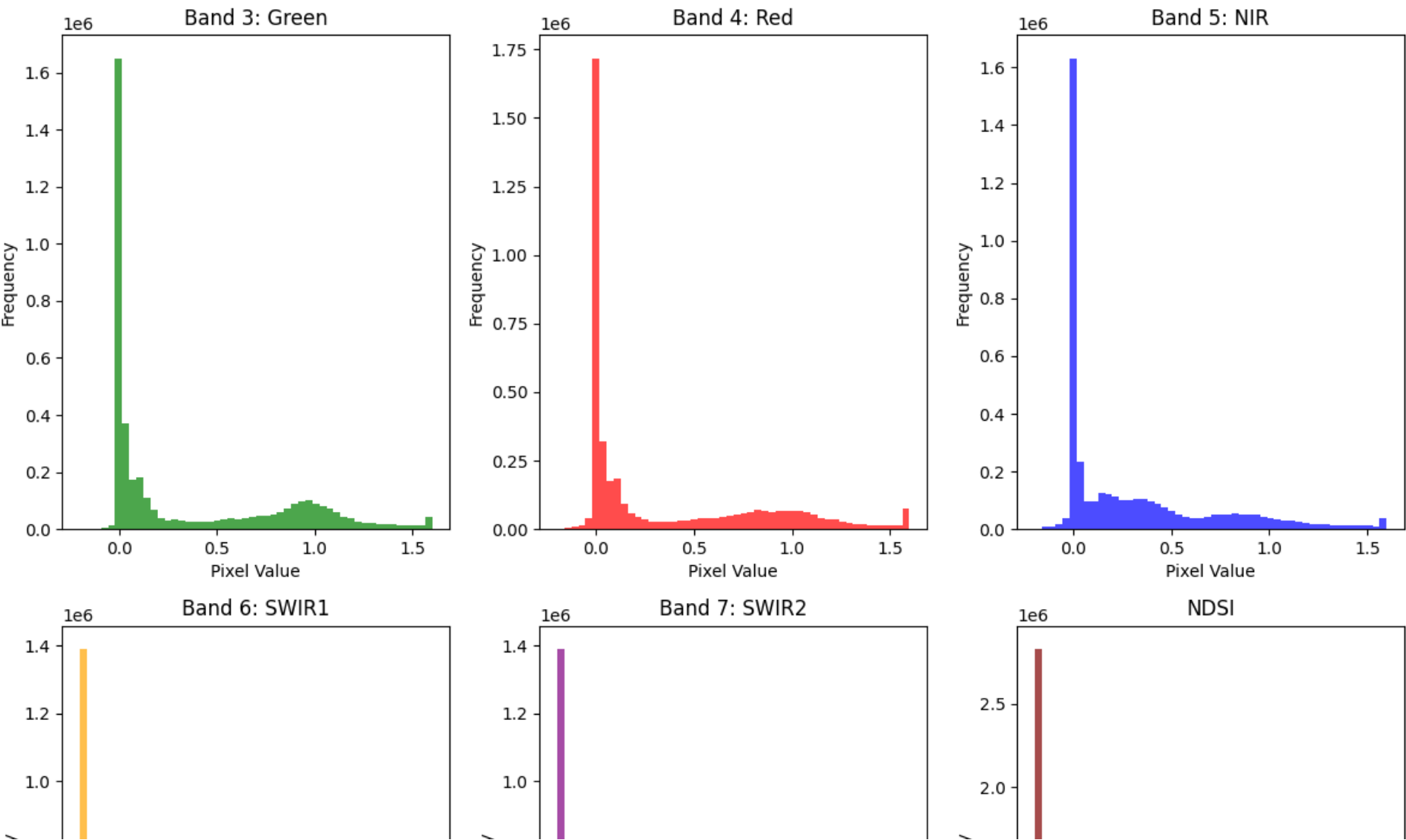


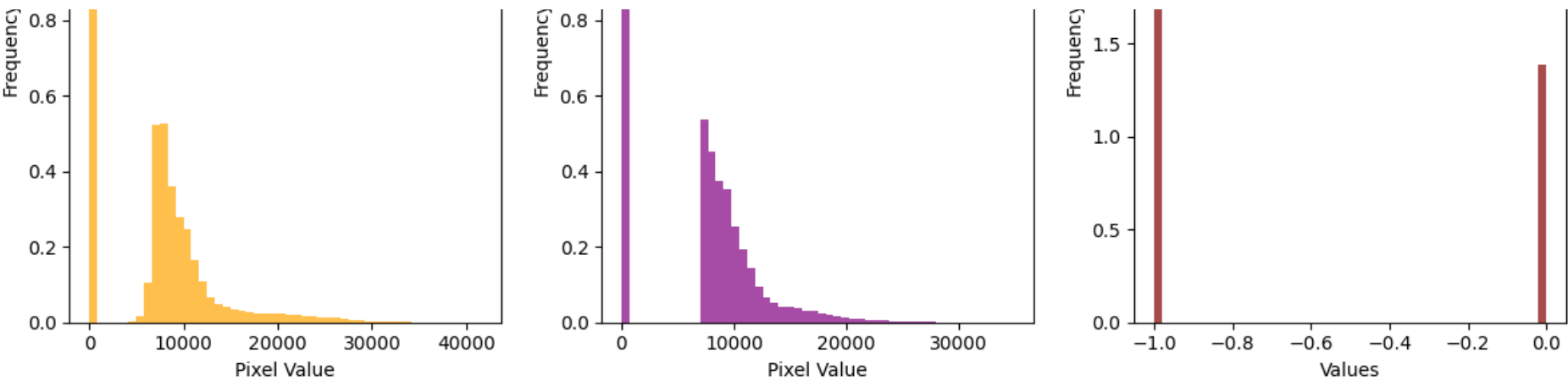
```

1 # Plotting the histogram to understand the data distribution, and what natural phenomenon is represented
2
3 # Summary of ALL BANDS and NDSI:
4 # Band 3 (Green) represents the visible green light reflected by vegetation, useful for vegetation analysis
5 # Band 4 (Red) corresponds to the visible red light and is sensitive to chlorophyll absorption, helpful for
6 # Band 5 (NIR) captures near-infrared light, with high reflectance from healthy vegetation, making it useful
7 # Band 6 (SWIR1) detects moisture content in soil, vegetation, and water, with lower reflectance from
8 # Band 7 (SWIR2) also measures moisture, with snow and ice showing low reflectance, and is used for
9 # The NDSI (Normalized Difference Snow Index) uses Band 3 and Band 6 to highlight snow and ice by comparing
10 # in the visible and SWIR bands; values close to 0 indicate non-snow surfaces, while higher values indicate
11
12 # here, normalizing any of the band values, can deteriorate the characteristics separation of glaciers
13
14 plt.figure(figsize=(12, 10))
15

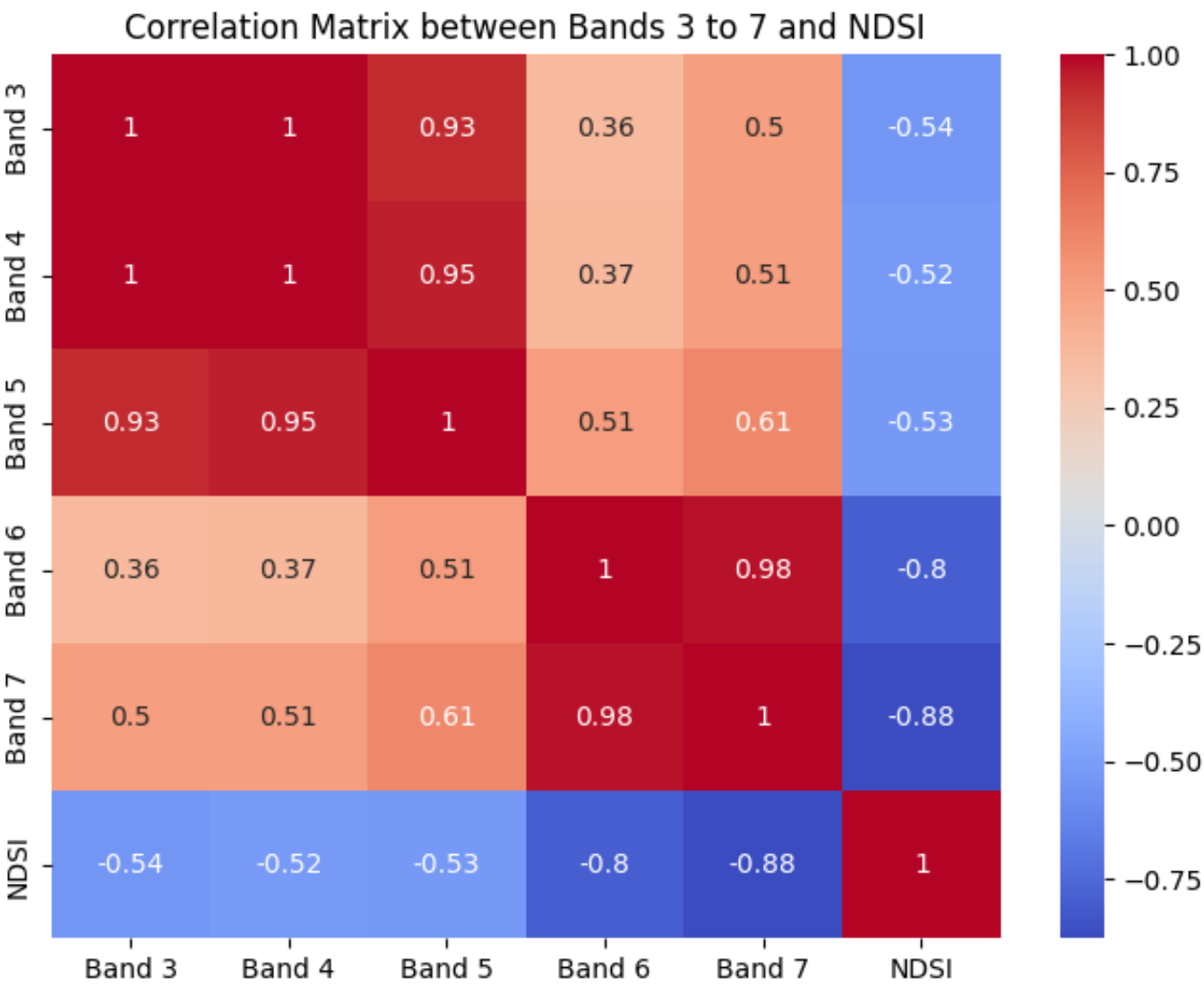
```

```
16 # Band 3: Green
17 plt.subplot(2, 3, 1)
18 plt.hist(band_3_flat, bins=50, color='green', alpha=0.7)
19 plt.title('Band 3: Green')
20 plt.xlabel('Pixel Value')
21 plt.ylabel('Frequency')
22
23 # Band 4: Red
24 plt.subplot(2, 3, 2)
25 plt.hist(band_4_flat, bins=50, color='red', alpha=0.7)
26 plt.title('Band 4: Red')
27 plt.xlabel('Pixel Value')
28 plt.ylabel('Frequency')
29
30 # Band 5: NIR
31 plt.subplot(2, 3, 3)
32 plt.hist(band_5_flat, bins=50, color='blue', alpha=0.7)
33 plt.title('Band 5: NIR')
34 plt.xlabel('Pixel Value')
35 plt.ylabel('Frequency')
36
37 # Band 6: SWIR1
38 plt.subplot(2, 3, 4)
39 plt.hist(band_6_flat, bins=50, color='orange', alpha=0.7)
40 plt.title('Band 6: SWIR1')
41 plt.xlabel('Pixel Value')
42 plt.ylabel('Frequency')
43
44 # Band 7: SWIR2
45 plt.subplot(2, 3, 5)
46 plt.hist(band_7_flat, bins=50, color='purple', alpha=0.7)
47 plt.title('Band 7: SWIR2')
48 plt.xlabel('Pixel Value')
49 plt.ylabel('Frequency')
50
51 plt.subplot(2, 3, 6)
52 plt.hist(NDSI_flat, bins=50, color='maroon', alpha=0.7)
53 plt.title('NDSI')
54 plt.xlabel('Values')
55 plt.ylabel('Frequency')
56
57 plt.tight_layout()
58 plt.show()
```





```
1 # Replotting the correlation matrix with NDSI included, and it is obvious that it will be highly cor
2 # But the high correlation between band 7 and NDSI can be due to the fact that there are similar spe
3 # SWIR (Shortwave Infrared) bands, particularly Band 6 (SWIR1) and Band 7 (SWIR2).
4 NDSI = np.nan_to_num(NDSI, nan=0)
5 NDSI_flat = NDSI.flatten()
6 data = np.stack([band_3_flat, band_4_flat, band_5_flat, band_6_flat, band_7_flat, NDSI_flat])
7
8 correlation_matrix = np.corrcoef(data)
9
10 # Create a heatmap of the correlation matrix
11 plt.figure(figsize=(8, 6))
12 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', xticklabels=["Band 3", "Band 4", "Band
13                               yticklabels=["Band 3", "Band 4", "Band 5", "Band 6", "Band 7", "NDSI"])
14 plt.title('Correlation Matrix between Bands 3 to 7 and NDSI')
15 plt.show()
```



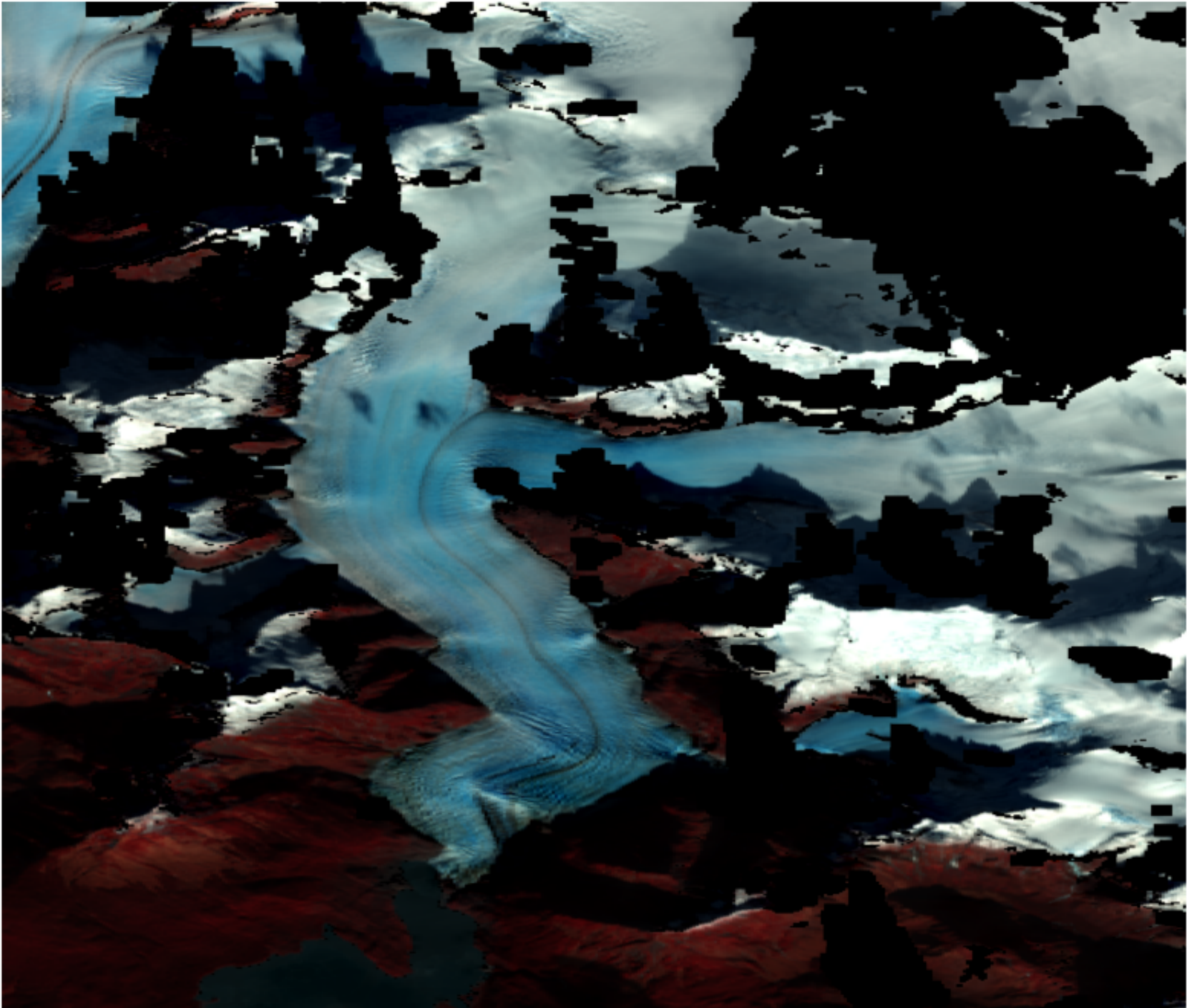

```
1 # Min value of NDSI being -1.0:
2 # This suggests that in some areas, the difference between the Green band (Band 3) and the SWIR band
3 # the reflectance in the SWIR band is much higher than in the Green band. The minimum value of -1.0
4 # extremely small or close to zero compared to the SWIR band in these areas.
5 # Max value of 0.0:
6 # The maximum value of 0.0 suggests that in all of the image, there are no areas where the Green re
7 # point that the result is positive.
8 # Mean of -0.67:
9 # A mean value of -0.67 indicates that, on average, the Green reflectance is significantly lower tha
10 # This is usually expected for non-snow/ice-covered areas (e.g., forests, urban areas), where the SW
11 # the Green reflectanc
12
13 NDSI = np.nan_to_num(NDSI, nan=0)
14 print("NDSI min:", NDSI.min(), "max:", NDSI.max(), "mean:", NDSI.mean())
```

NDSI min: -1.0000713605834395 max: 0.0 mean: -0.6704336916204889

```
1 # Using the three bands to plot an RGB image of glacier characteristics measured over 2 months
2 band_5_norm = band_5 / band_5.max()
3 band_4_norm = band_4 / band_4.max()
4 band_3_norm = band_3 / band_3.max()
5
6 rgb_image = np.stack((band_5_norm, band_4_norm, band_3_norm), axis=-1)
7
8 plt.figure(figsize=(10, 10))
9 plt.imshow(rgb_image)
10 plt.title("False Color Composite (NIR, Red, Green)")
11 plt.axis('off')
12 plt.show()
13
```

WARNING:matplotlib.image:Clipping input data to the valid range for imshow with RGB data ([0..1] for float)

False Color Composite (NIR, Red, Green)



1 # Further Steps: I will either use the NDSI image or the RGB Image to perform segmentation (Model p
2 # (may be with Swin Backbone)), and then I will perform classification, on initial stage I will keep
3 # (whether snow is present or absent), which if time permits will further include classes like water