# B

# Code of Project

All code implemented is publicly accessible through a Git-lab repository located at following link: http

**Listing B.1:** Code for NCSS: Creation of table "ids_loc" through the "sqlite3" command using the machine's terminal. Second command counts the number of rows

```sql
1  CREATE TABLE ids_loc AS
2  SELECT *
3  FROM lab_layer
4  JOIN lab_site ON lab_layer.site_key = lab_site.site_key;
5
6  --Counting ids_loc rows
7  SELECT COUNT(*) FROM ids_loc;
```

**Listing B.2:** Code for NCSS: Creation of table "ids_loc_prop" through the "sqlite3" command using the machine's terminal. Second command counts the number of rows.

```
1  --Appending lab_chemical_properties table
2  CREATE TABLE ids_loc_chem AS
3  SELECT *
4  FROM ids_loc
5  JOIN lab_chemical_properties
6  ON ids_loc.labsampnum = lab_chemical_properties.labsampnum
7  AND ids_loc.layer_key = lab_chemical_properties.layer_key;
8
9  --Appending lab_physical_properties table
10 CREATE TABLE ids_loc_chem_phy AS
11 SELECT *
12 FROM ids_loc_chem
13 JOIN lab_physical_properties
14 ON ids_loc_chem.labsampnum = lab_physical_properties.labsampnum
15 AND ids_loc_chem.layer_key = lab_physical_properties.layer_key;
16
17 --Appending lab_major_and_trace_elements_and_oxides table
18 CREATE TABLE ids_loc_prop AS
19 SELECT *
20 FROM ids_loc_chem_phy
21 JOIN lab_major_and_trace_elements_and_oxides
22 ON ids_loc_chem_phy.labsampnum = lab_major_and_trace_elements_and_oxides.labsampnum
23 AND ids_loc_chem_phy.layer_key = lab_major_and_trace_elements_and_oxides.layer_key;
24
25 --Counting ids_loc_prop rows
26 SELECT COUNT(*) FROM ids_loc_prop;
```

**Listing B.3:** Code for NCSS: Creation of table "ids_loc_prop_time" through the "sqlite3" command using the machine's terminal. Second command counts the number of rows

```
1  CREATE TABLE ids_loc_prop_time AS
2  SELECT *
3  FROM ids_loc_prop
4  JOIN lab_pedon ON ids_loc_prop.pedon_key = lab_pedon.pedon_key
5  AND ids_loc_prop.site_key = lab_pedon.site_key;
6
7  --Counting ids_loc rows
8  SELECT COUNT(*) FROM ids_loc_prop_time;
```

**Listing B.4:** Code for NCSS: Creation of table "data" through the "sqlite3" command using the machine's terminal.

```sql
CREATE TABLE data AS
SELECT labsampnum, pedlabsampnum, hzn_top, hzn_bot, observation_date,
latitude_std_decimal_degrees, longitude_std_decimal_degrees, total_carbon_ncs,
total_nitrogen_ncs, phosphorus_bray1, phosphorus_bray2, phosphorus_major_element,
phosphorus_trace_element, k_nh4_ph_7, potassium_major_element, ca_nh4_ph_7,
calcium_major_element, mg_nh4_ph_7, magnesium_major_element, total_sulfur_ncs,
copper_trace_element, fe_ammoniumoxalate_extractable, iron_sodium_pyro_phosphate,
iron_major_element, manganese_ammonium_oxalate, manganese_dithionite_citrate,
manganese_kcl_extractable, manganese_major_element, manganese_trace_element,
molybdenum_trace_element, zin_trace_element
FROM ids_loc_prop_time;
```

**Listing B.5:** Code for NCSS: dropping rows with missing latitude, longitude or date

```python
import pandas as pd
# DataFrame with NCSS data
df = pd.read_csv('./NCSS/data.csv')
df = df.dropna(subset=['observation_date', 'latitude_std_decimal_degrees', '
    longitude_std_decimal_degrees'])
```

**Listing B.6:** Code for NCSS: depth parsing

```python
df = df.dropna(how='all', subset=['hzn_top', 'hzn_bot'])
# Filter rows where more than 50% of the depth interval is below 30
df = df[(30 - df['hzn_top']) / (df['hzn_bot'] - df['hzn_top']) >= 0.5]
# Drop duplicates based on 'pedlabsampnum' keeping the one with the lowest '
    hzn_top'
df = df.sort_values(by='hzn_top').drop_duplicates(subset='pedlabsampnum',
    keep='first')
```

**Listing B.7:** Code for NCSS: nutrient conversion to ppm and averaging

```python
# Features that need gravimetric to ppm convertion
gravimetric_features = (
    'total_carbon_ncs',
```

```python
4        'total_nitrogen_ncs',
5        'total_sulfur_ncs',
6        'fe_ammoniumoxalate_extractable',
7        'iron_sodium_pyro_phosphate',
8        'manganese_dithionite_citrate'
9        )
10
11  # Features that need meq/100g to ppm convertion
12  meq_per_100_features = {
13       'k_nh4_ph_7': 39.0,
14       'ca_nh4_ph_7': 40.08,
15       'mg_nh4_ph_7': 24.305
16  }
17
18  # Applying convertions according to README.md file
19  # Convert all features in gravimetric percentage
20  for feature in gravimetric_features:
21       df[feature] = df[feature].apply(lambda x: round(x * 10**4, 3))
22  # Convert all features in meq/100g
23  for feature in meq_per_100_features:
24       df[feature] = df[feature].apply(lambda x: round(x * 10 * meq_per_100
            _features[feature], 3))
25
26  # Compute average between features that tell the same info. nan values are
         removed from equation
27  df['P'] = df[['phosphorus_bray1', 'phosphorus_bray2', '
         phosphorus_major_element', 'phosphorus_trace_element']].mean(axis=1)
28  df['K'] = df[['k_nh4_ph_7', 'potassium_major_element']].mean(axis=1).round(3)
29  df['Ca'] = df[['ca_nh4_ph_7', 'calcium_major_element']].mean(axis=1).round(3)
30  df['Mg'] = df[['mg_nh4_ph_7', 'magnesium_major_element']].mean(axis=1).round(
         3)
31  df['Fe'] = df[['fe_ammoniumoxalate_extractable', 'iron_sodium_pyro_phosphate'
         , 'iron_major_element']].mean(axis=1).round(3)
32  df['Mn'] = df[['manganese_ammonium_oxalate', 'manganese_dithionite_citrate',
         'manganese_kcl_extractable', 'manganese_major_element', '
         manganese_trace_element']].mean(axis=1).round(3)
33  df['Cec'] = df[['cec7_clay_ratio', 'cec_nh4_ph_7']].mean(axis=1).round(3)
34
```

```
35  # Renaming
36  df = df.rename(columns={
37      'observation_date': 'date',
38      'latitude_std_decimal_degrees': 'latitude',
39      'longitude_std_decimal_degrees': 'longitude',
40      'total_carbon_ncs': 'C',
41      'total_nitrogen_ncs': 'N',
42      'total_sulfur_ncs': 'S',
43      'copper_trace_element': 'Cu',
44      'molybdenum_trace_element': 'Mo',
45      'zinc_trace_element': 'Zn'
46  })
47
48  # Keeping necessary features
49  df = df[[
50      'labsampnum', 'pedlabsampnum', 'date',
51      'latitude', 'longitude',
52      'C', 'N', 'P', 'K', 'Ca', 'Mg', 'S', 'Cu', 'Fe', 'Mn', 'Mo', 'Zn','Cec'
53      ]]
```

**Listing B.8:** Code for NCSS: nutrient missing values and percentage

```
1  nutrients = ['C', 'N', 'P', 'K', 'Ca', 'Mg', 'S', 'Cu', 'Fe', 'Mn', 'Mo', 'Zn
      ']
2  # Calculate missing value count per feature
3  missing_count_per_feature = df[nutrients].isnull().sum()
4  # Calculate missing value percentage per feature
5  missing_percentage_per_feature = round((missing_count_per_feature / len(df))
      * 100, 2)
6  # Create a DataFrame to display the results
7  pd.DataFrame({
8      'Missing Count': missing_count_per_feature,
9      'Missing Percentage': missing_percentage_per_feature
10 }).transpose()
```

**Listing B.9:** Code for NCSS: removing points with no nutrient data

```
1 df = df.dropna(how='all', subset=nutrients)
```

**Listing B.10:** Code for NCSS: count data for each satellite's life span

```
1  df['date'] = pd.to_datetime(df['date'], format='%Y-%m-%d').dt.date
2
3  mask1 = (df['date'] >= datetime.date(1972, 7, 1)) & (df['date'] <= datetime.
       date(1992, 10, 31))
4  mask2 = (df['date'] >= datetime.date(2012, 6, 1)) & (df['date'] <= datetime.
       date(2013, 1, 31))
5
6  time_interval = [
7      mask1 | mask2,
8      (df['date'] >= datetime.date(1982, 7, 1)) & (df['date'] <= datetime.date(
           2012, 5, 31)),
9      (df['date'] >= datetime.date(1999, 4, 1)),
10     (df['date'] >= datetime.date(2013, 2, 1)),
11     (df['date'] >= datetime.date(2017, 1, 1))
12     ]
13
14 MSS = df[time_interval[0]]
15 TM = df[time_interval[1]]
16 ETM = df[time_interval[2]]
17 OLI = df[time_interval[3]]
18 S2 = df[time_interval[4]]
```

**Listing B.11:** Code for NCSS: EOS extraction from LMSS, LTM and LETM

```
1  import terrasensetk as tstk
2  import sentinelhub as sh
3
4  # Sentinel hub credentials
5  config = sh.SHConfig()
6  # Id of the client
7  config.sh_client_id = "------------------------"
8  # Sectret token for client
9  config.sh_client_secret = "---------------------"
```

```
10  config.save ()
11
12  # Reading shapefile
13  shp_path = '../NCSS/data.shp'
14  nutrients = ['C', 'N', 'P', 'K', 'Ca', 'Mg', 'S', 'Cu', 'Fe', 'Mn', 'Mo', 'Zn
       ']
15  # Landsat 1-5 MSS
16  mss_down = tstk.LandsatMSSDownloader(config=config, shapefile=shp_path,
       targets=nutrients)
17  # Landsat 4-5 TM
18  tm_down = tstk.LandsatTMDownloader(config=config, shapefile=shp_path, targets
       =nutrients)
19  # Landsat 7 ETM+
20  etm_down = tstk.LandsatETMDownloader(config=config, shapefile=shp_path,
       targets=nutrients)
21
22  downloaders = [mss_down, tm_down, etm_down]
23  for d in downloaders:
24      d.download(path=f'./NCSS/EOPatches/{d.name}',
25              date_fmt='%d/%m/%Y', padding_days=10, maxcc=0.1, res=20,
                  report=True)
```

**Listing B.12:** Code for LUCAS: EOS extraction from S2

```
1  import terrasensetk as tstk
2  import sentinelhub as sh
3
4  # Sentinel hub credentials
5  config = sh.SHConfig()
6  # Id of the client
7  config.sh_client_id = "-----------------------"
8  # Sectret token for client
9  config.sh_client_secret = "---------------------"
10 config.save ()
11
12 # Reading shapefile
13 shp_path = '../LUCAS/data.shp'
14 nutrients = ['N', 'P', 'K']
```

```
15  # Sentinel S2-A and S2-B
16  s2_down = tstk.Sentinel2Downloader(config=config, shapefile=shp_path, targets
        =nutrients)
17  # Downloading
18  s2_down.download(path=f'./NCSS/EOPatches/{d.name}',
19          date_fmt='%d/%m/%Y', padding_days=10, maxcc=0.1, res=20, report=True)
```

**Listing B.13:** Code for LUCAS: count and percentage number of missing nutrients

```
1   import pandas as pd
2
3   eopatches_path = '../LUCAS/EOPatches/sentinel-2'
4   # Creates dataset from EOPatches
5   dataset = tstk.Dataset(eopatches_path)
6
7   df = dataset.get_df()
8   # Forced convertion to float. < LOD string beomces nan
9   df.loc[:, labels] = df[labels].apply(lambda x: pd.to_numeric(x, errors='
        coerce'))
10  df[labels] = df[labels].astype(float)
11
12  pd.DataFrame(data=(df[labels].count(), df.shape[0] - df[labels].count()))
```

**Listing B.14:** Code for LUCAS: scatterplot for N, P and K

```
1   import matplotlib.pyplot as plt
2
3   for l, name in zip(labels, ('Nitrogen', 'Phosphorus', 'Potassium')):
4       values = df[l]
5       plt.scatter(range(values.shape[0]), values, label=name)
6       plt.title(f'Distribution of {name} ({l})')
7       plt.xlabel('Index')
8       plt.ylabel(f'{l} (ppm)')
9       plt.show()
```

**Listing B.15:** Code for LUCAS: nutrients outlier count

```python
data = []

for nut in labels:
    Q1 = df[nut].quantile(0.25)
    Q3 = df[nut].quantile(0.75)
    # Calculate interquartile range (IQR)
    IQR = Q3 - Q1
    # Define lower and upper bounds for outliers
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    # Count outliers
    data.append(df[nut][(df[nut] < lower_bound) | (df[nut] > upper_bound)].
        count())

display(pd.DataFrame(data=data, columns=['# Outliers (IQR)'], index=labels))
```

**Listing B.16:** Code for LUCAS: histogram of "LC0_Desc" classes

```python
# Gathering unique count data
df_data = df['LC0_Desc'].value_counts().reset_index()

plt.figure(figsize=(8, 8))
plot = sns.barplot(x="LC0_Desc", y="count", data=df_data)

for bar in plot.patches:
    # x-coordinate: bar.get_x() + bar.get_width() / 2
    # y-coordinate: bar.get_height()
    # free space to be left to make graph pleasing: (0, 8)
    # ha and va stand for the horizontal and vertical alignment
    plot.annotate(format(bar.get_height(), '.2f'),
                  (bar.get_x() + bar.get_width() / 2, bar.get_height()),
                  ha='center', va='center', size=15, xytext=(0, 8),
                     textcoords='offset points')

plt.xlabel("Land Cover 0", size=14)
plt.ylabel("Count", size=14)
```

```
18 plt.title("Land Cover Class frequency")
19 plt.show()
```

**Listing B.17:** Code for LUCAS: histogram of "LC1_Desc" classes

```
1 # Gathering unique count data
2 df_data = df['LC1_Desc'].value_counts().reset_index()
3
4 # Defining the plot size
5 plt.figure(figsize=(8, 8))
6 # Defining the values for x-axis, y-axis
7 # and from which dataframe the values are to be picked
8 plot = sns.barplot(y="LC1_Desc", x="count", data=df_data)
9 plt.xlabel("Count", size=14)
10 plt.ylabel("Land Cover 1", size=14)
11 plt.title("Land Cover 1 frequency")
12 plt.show()
```

**Listing B.18:** Code for LUCAS: table with N, P and K data summary and their outliers according to the IQR method

```
1 data = []
2
3 # For each nutrient
4 for nut in labels:
5     # Calculate quartiles
6     Q1 = df[nut].quantile(0.25)
7     Q3 = df[nut].quantile(0.75)
8     # Calculate interquartile range (IQR)
9     IQR = Q3 - Q1
10    # Define lower and upper bounds for outliers
11    lower_bound = Q1 - 1.5 * IQR
12    upper_bound = Q3 + 1.5 * IQR
13    # Count outliers
14    data.append(df[nut][(df[nut] < lower_bound) | (df[nut] > upper_bound)].
          count())
15
16 display(pd.DataFrame(data=data, columns=['# Outliers (IQR)'], index=labels))
```

```
17
18 df[labels].describe()
```

**Listing B.19:** Code for final dataset: Conversion to classification

```python
1  # Fertility limits
2  N_fertility_limits = [10, 50]
3  P_fertility_limits = [10.9, 21.4]
4  K_fertility_limits = [40, 80]
5
6  # Auxiliar function for classification
7  def classify(x, lower, upper):
8      if lower <= x <= upper:
9          return 1.0
10     elif x < lower or x > upper:
11         return 0.0
12     return x
13
14 # Classify each nutrient with function
15 for l, limits in zip(labels, (N_fertility_limits, P_fertility_limits,
       K_fertility_limits)):
16     df[l] = df[l].apply(classify, args=limits)
```

**Listing B.20:** Code for final dataset: Balance assessment

```python
1  data = []
2  N_total, N_fert, N_inf = df['N'].count(), df['N'][df['N'] == 0.0].count(), df
       ['N'][df['N'] == 1.0].count()
3  P_total, P_fert, P_inf = df['P'].count(), df['P'][df['P'] == 0.0].count(), df
       ['P'][df['P'] == 1.0].count()
4  K_total, K_fert, K_inf = df['K'].count(), df['K'][df['K'] == 0.0].count(), df
       ['K'][df['K'] == 1.0].count()
5
6  data.append((N_total, N_fert, N_fert/N_total, N_inf, N_inf/N_total))
7  data.append((P_total, P_fert, P_fert/P_total, P_inf, P_inf/P_total))
8  data.append((K_total, K_fert, K_fert/K_total, K_inf, K_inf/K_total))
9
```

```
10 pd.DataFrame(data=data, columns=('X_l', 'infertile', 'infertile (%)', '
       fertile', 'fertile (%)'), index=labels)
```

**Listing B.21:** Code for final dataset: Balance scatterplot

```
1 for l, name in zip(labels, ('Nitrogen', 'Phosphorus', 'Potassium')):
2     values = df[l]
3     if l == 'N':
4         lower, upper = N_fertility_limits
5     elif l == 'P':
6         lower, upper = P_fertility_limits
7     if l == 'K':
8         lower, upper = K_fertility_limits
9     # Define colors based on conditions
10    # Define colors based on conditions
11    colors = np.where((values < lower) | (values > upper), 'red', 'green')
12    # Plot infertile points in red
13    plt.scatter(np.where(colors == 'red')[0], values[colors == 'red'], color=
          'red', label=f'Infertile {name}')
14    # Plot fertile points in green
15    plt.scatter(np.where(colors == 'green')[0], values[colors == 'green'],
          color='green', label=f'Fertile {name}')
16
17    plt.title(f'Distribution of {name} ({l})')
18    plt.xlabel('Index')
19    plt.ylabel(f'{l} (mg/kg)')
20    # Add a legend and label for color interpretation
21    plt.legend()
22    plt.show()
```

**Listing B.22:** Code for final dataset: Trinary classification

```
1 # Auxiliar function for classification
2 def classify(x, lower, upper):
3     if lower <= x <= upper:
4         return 'sufficient'
5     if x < lower:
```

```
6          return 'defice'
7      if x > upper:
8          return 'toxic'
9      return x
10
11 # Classify each nutrient with aux function
12 for l, limits in zip(labels, (N_fertility_limits, P_fertility_limits,
       K_fertility_limits)):
13     df[l] = df[l].apply(classify, args=limits)
```

**Listing B.23:** Code for final dataset: Balance scatterplot for trinary P and K

```
1 for l, name in zip(('P', 'K'), ('Phosphorus', 'Potassium')):
2      values = df[l]
3      if l == 'P':
4          lower, upper = P_fertility_limits
5      elif l == 'K':
6          lower, upper = K_fertility_limits
7      # Define colors based on conditions
8      defice = np.where(values < lower)
9      fertile = np.where((values >= lower) & (values <= upper))
10     toxic = np.where(values > upper)
11
12     # Plot defice points in yellow
13     plt.scatter(defice[0], values[values < lower], color='yellow', label=f'
           Defice')
14     # Plot defice points in yellow
15     plt.scatter(fertile[0], values[(values >= lower) & (values <= upper)],
           color='green', label=f'Fetile')
16     # Plot defice points in yellow
17     plt.scatter(toxic[0], values[values > upper], color='black', label=f'
           Toxic')
18
19     plt.title(f'Distribution of {name} ({l})')
20     plt.xlabel('Index')
21     plt.ylabel(f'{l} (mg/kg)')
22     # Add a legend and label for color interpretation
23     plt.legend()
```

```
24    plt.show()
```

**Listing B.24:** Code for final dataset: Balance table for trinary P and K

```python
1  data = []
2
3  N_total, N_insuf, N_fert, N_tox = df['N'].count(), df['N'][df['N'] == '
       sufficient'].count(), df['N'][df['N'] == 'defice'].count()\
4      , df['N'][df['N'] == 'toxic'].count()
5  P_total, P_insuf, P_fert, P_tox = df['P'].count(), df['P'][df['P'] == '
       sufficient'].count(), df['P'][df['P'] == 'defice'].count()\
6      , df['P'][df['P'] == 'toxic'].count()
7  K_total, K_insuf, K_fert, K_tox = df['K'].count(), df['K'][df['K'] == '
       sufficient'].count(), df['K'][df['K'] == 'defice'].count()\
8      , df['K'][df['K'] == 'toxic'].count()
9
10 data.append((N_total, N_insuf, N_insuf/N_total, N_fert, N_fert/N_total, N_tox
       , N_tox/N_total))
11 data.append((P_total, P_insuf, P_insuf/P_total, P_fert, P_fert/P_total, P_tox
       , P_tox/P_total))
12 data.append((K_total, K_insuf, K_insuf/K_total, K_fert, K_fert/K_total, K_tox
       , K_tox/K_total))
13
14 pd.DataFrame(data=data, columns=('X_l',
15                                   'defice', 'defice (%)',
16                                   'fertile', 'fertile (%)',
17                                   'toxic', 'toxic (%)'), index=labels)
```

**Listing B.25:** Prior code for tstk: get_bbox() and get_bbox_with_data()

```python
1
2      def get_bbox_with_data(self):
3          """
4
5          Returns:
6              GeoDataFrame: Contains the bboxes which have associated
                   groundtruth
```

```python
7          """
8          if self._bbox_with_groundtruth is not None:
9              return self._bbox_with_groundtruth
10         self._bbox_with_groundtruth = self.get_groundtruth().copy(deep=True).
               reset_index()
11         self._bbox_with_groundtruth.geometry = self.get_bbox().reset_index().
               geometry

12

13         return self._bbox_with_groundtruth

14

15

16     def get_bbox(self,buffer=0.005,reset=False):

17

18         """
19         Creates a grid of bbox over the dataset

20

21         Args:
22             dataset (GeoDataFrame): [description]
23             expected_bbox_size (int, optional): The desired size of the bbox
                   in meters. Defaults to 2000.
24             reset (bool, optional): Wether it should recalculate the
                   bbox_list. Defaults to False.

25

26         Returns:
27             GeoDataFrame: GeoDataFrame of the dataset divided in square bbox
                   of size of expected_bbox_size.
28         """

29

30         #create bboxes around the groundtruth
31         if self._dataset_bbox is not None and not reset:
32             return self._dataset_bbox
33         points_of_interest = [shapely.geometry.MultiPolygon([i.centroid.
               buffer(0.00001) for i in self.get_groundtruth().geometry.values])
               ]
34         points_grid = self.get_groundtruth().geometry.apply(lambda x: BBox(x.
               centroid.buffer(buffer), sh.CRS.WGS84)).to_list()
35         self.dataset_bbox_splitter = CustomGridSplitter(points_of_interest,
36             sh.CRS.WGS84.pyproj_crs(),
```

```
37              points_grid )
38          geometry = [ Polygon ( bbox . get_polygon () ) for bbox in self .
                dataset_bbox_splitter . get_bbox_list () ]
39          self . _dataset_bbox = gpd . GeoDataFrame ( crs = sh . CRS . WGS84 . pyproj_crs () ,
                geometry = geometry )
40          #self._dataset_bbox = self._dataset_bbox.drop_duplicates()
41          return self . _dataset_bbox
```

**Listing B.26:** Code for tstk: Downloader constructure

```
1 # Creates dataset from shapefile
2 if shapefile is not None :
3      self . dataset = gpd . read_file ( shapefile )
4      if groundtruth_col is not None :
5          # Geometry keeps groundtruth geometry values
6          has_geom = ~ self . dataset [ groundtruth_col ]. isna ()
7          self . dataset . loc [ has_geom , 'geometry' ] = loads ( self . dataset . loc [
                has_geom , groundtruth_col ])
8          self . dataset . drop ( columns = groundtruth_col , inplace = True )
```

**Listing B.27:** Code for tstk: Function to_square in utils.util.py

```
1 def to_square ( center_point : Point , area : int | float ) -> Polygon :
2      """
3      #### Function that returns a Polygon of a square with center in Point \
4      and with area in meters^2.
5
6      Args :
7          center_point : Point to use as the center of the square .
8          area : Area in m2 .
9
10     Returns :
11         Polygon with all corner points of square
12
13     Raises ValueError and TypeError
14     """
15     if not isinstance ( center_point , Point ):
```

```
16        raise TypeError(f'point must be Point. Got {type(center_point)}.')
17    if type(area) is not int and type(area) is not float:
18        raise TypeError(f'area must be int or float. Got {type(area)}.')
19
20    l = sqrt(area)
21    return center_point.buffer(l, cap_style='square')
```

**Listing B.28:** Code for tstk: Downloader.calculate_area()

```
1 def calculate_area(self, area: float|int = 2000) -> gpd.GeoDataFrame:
2        """
3        #### Function that substitutes each Point found in geometry\
4        with a square Polygon with geom = 'area' and center in Point.
5
6        Args:
7            area: int or float equal to the area to be calculated.
8
9        Return:
10            GeoDataFrame
11
12        Raises TypeError or ValueError
13        """
14        if type(area) is not int and type(area) is not float:
15            raise TypeError(f'Area must be int or float. Got {type(area)}.')
16        if area <= 0:
17            raise ValueError(f'Area must be > 0. Got {area}.')
18        self.dataset.to_crs('EPSG:3857', inplace=True)
19        # apply function to all Point
20        self.dataset['geometry'] = self.dataset['geometry'].apply(lambda g:
            to_square(g, area) if isinstance(g, Point) else g
21        )
22        self.dataset.to_crs(CRS.WGS84.pyproj_crs(), inplace=True)
23        return self.dataset
```

**Listing B.29:** Code for tstk: Downloader.show_geometries()

```
1 def show_geometries(self) -> Map:
```

```
2          """
3          #### Function that returns global map with all geometries
4
5          Returns:
6              folium.Map
7          """
8          loc = self.dataset.unary_union.centroid.y, self.dataset.unary_union.
               centroid.x
9          mapa = Map(location=loc, zoom_start=3)
10         for g in self.dataset['geometry']:
11             GeoJson(mapping(g)).add_to(mapa)
12         return mapa
```

**Listing B.30:** Prior code for tstk: Downloader.download_images()

```
1 def download_images(self,path,bands=None,subset=None,date_field="SURVEY_DATE"
       ,date_fmt='%d/%m/%y',padding_days=2,maxcc=0.5):
2          """Downloads the specified images into the users filesystem.
3
4          Args:
5              path (str): Path to where the dataset should be saved
6              subset (DataFrame, optional): Slice of the dataframe returned by
                  `get_bbox_with_data()`.
7          """
8          if bands is None:
9              bands = ["B01","B02","B03","B04","B05","B06","B07","B08","B8A","
                  B09","B11","B12"]
10         if subset is None:
11             subset = self.get_bbox_with_data()
12         if not os.path.isdir(path):
13             os.makedirs(path)
14
15         add_data = SentinelHubInputTask(
16             bands_feature=(FeatureType.DATA, 'BANDS'),
17             resolution=10,
18             bands=bands,
19             maxcc=maxcc,
20             time_difference=datetime.timedelta(minutes=120),
```

```python
21              data_collection=DataCollection.SENTINEL2_L2A,
22              additional_data=[(FeatureType.MASK, 'dataMask'),
23                       (FeatureType.MASK, 'CLM'),
24                       (FeatureType.DATA, 'CLP')],
25              max_threads=5
26          )
27          add_data_node = EONode(add_data,name="add_data_node")
28
29          add_vector = AddFeatureTask((FeatureType.VECTOR_TIMELESS,"LOCATION"))
30          add_vector_node = EONode(add_vector,[add_data_node],name="
                add_vector_node")
31          #add_lucas = AddFeature((FeatureType.META_INFO,"LUCAS_DATA"))
32          #to get the surrounding data, one can apply a buffered vector to
                raster and set the non overlapped value some value to distinguish
33          add_raster_buffer = VectorToRasterTask((FeatureType.VECTOR_TIMELESS,"
                LOCATION"),(FeatureType.MASK_TIMELESS,"IS_VALID"), values = 5,
                buffer=0.0005, raster_shape=(FeatureType.MASK, 'CLM'),
                no_data_value=0,raster_dtype=np.uint8)
34          add_raster_buffer_node = EONode(add_raster_buffer,[add_vector_node],
                name="add_raster_buffer_node")
35          add_raster = VectorToRasterTask((FeatureType.VECTOR_TIMELESS,"
                LOCATION"),(FeatureType.MASK_TIMELESS,"IS_VALID"), values = 1,
                raster_shape=(FeatureType.MASK, 'CLM'),write_to_existing = True,
                no_data_value=0,raster_dtype=np.uint8)
36          add_raster_node = EONode(add_raster,[add_raster_buffer_node],name="
                add_raster_node")
37
38          concatenate = MergeFeatureTask({FeatureType.DATA: ['BANDS']},(
                FeatureType.DATA, 'FEATURES'))
39          concatenate_node = EONode(concatenate,[add_raster_node],name="
                concatenate_node")
40          save = SaveTask(path, overwrite_permission=OverwritePermission.
                OVERWRITE_PATCH)
41          save_node = EONode(save,[concatenate_node],name="save_node")
42          workflow = EOWorkflow([add_data_node,add_vector_node,
                add_raster_buffer_node,add_raster_node,concatenate_node,save_node
                ])
43
```

```
44          execution_args = []
45          for id, wrap_bbox in enumerate(subset.iterrows()):
46              i, bbox = wrap_bbox
47
48              time_interval = (get_time_interval(bbox[date_field],padding_days,
                    date_fmt=date_fmt))
49              gdf = gpd.GeoDataFrame(bbox)
50              gdf = gpd.GeoDataFrame(gdf.transpose())
51              gdf = gdf.rename(columns={0:'geometry'}).set_geometry('geometry')
52              # gdf.set_geometry('geometry')
53              gdf.crs = sh.CRS.WGS84.pyproj_crs()
54
55              lucas_points_intersection = self.get_groundtruth()[self.
                    get_groundtruth().geometry.values.intersects(gdf.geometry.
                    values[0])]
56              execution_args.append({
57                  add_vector_node:{'data': lucas_points_intersection},
58                  add_data_node:{'bbox': BBox(bbox.geometry, crs=self.dataset.
                        crs), 'time_interval': time_interval},
59                  save_node: {'eopatch_folder': f'eopatch_{id}'}
60              })
61          executor = EOExecutor(workflow, execution_args, save_logs=True)
62          executor.run(workers=5, multiprocess=False)
63
64          executor.make_report()
```

**Listing B.31:** Code for tstk: Downloader.download()

```
1 def download(self, path: str, subset: gpd.GeoDataFrame = None,
2               date_fmt: str = '%d/%m/%y', padding_days: int = 2, maxcc:
                    float = 0.2, res: int = 10, report: bool = False, bands:
                    List[str] = None):
3       """
4       #### Downloads bands to EOPatches in 'path'.
5
6       Args:
7           path: Path to where the data is saved
8           subset: Slice of the dataframe. Default is all of it
```

```python
 9              date_frm: String representing the format of the date stored in df
10              padding_days: Number of days subtracted and added to create the
                    time interval
11              maxcc: Percentage from 0-1 of cloud coverage allowed in data
                    extracted
12              res: Resolution of images. I.e. meters/pixel
13              report: When true, saves the report as an html file
14              bands: List of bands that user intends to extract. None is all
                    available

16          NOTE:
17              - Bands are specified according to sentinel hub examples.
18              - Entries where groundtruth_col is Nan will be computed according
                    to are.
19          """
20          if type(path) is not str:
21              raise TypeError(f'path must be str. Got {type(path)}')
22          if subset is not None and not isinstance(subset, gpd.GeoDataFrame):
23              raise TypeError(f'subset must be geopandas.GeoDataFrame. Got {
                    type(subset)}')
24          else:
25              subset = self.dataset
26          if date_fmt is not None and type(date_fmt) is not str:
27              raise TypeError(f'date_fmt must be str. Got {type(date_fmt)}')
28          if padding_days is not None and type(padding_days) is not int:
29              raise TypeError(f'padding_days must be int. Got {type(
                    padding_days)}')
30          if maxcc is not None and type(maxcc) is not float:
31              raise TypeError(f'maxcc must be float. Got {type(maxcc)}')
32          if res is not None and type(res) is not int:
33              raise TypeError(f'res must be int. Got {type(res)}')
34          if report is not None and type(report) is not bool:
35              raise TypeError(f'report must be bool. Got {type(report)}')
36          if bands is not None and type(bands) is not list:
37              raise TypeError(f'bands must be list. Got {type(bands)}')

39          # Parses requested bands
40          if bands is None:
```

```python
41          bands = list(self.bands.values())
42      else:
43          for band in bands:
44              if type(band) is not str:
45                  raise TypeError(f'band must be str. Got {type(band)}')
46          if not set(bands).issubset(self.bands.keys()):
47              raise ValueError(f'Only {self.bands.keys()} bands are
                    available. Got {bands}')
48          bands = [b for b in self.bands.values() if b in bands]
49
50      # Creates dir if not existing
51      if not os.path.isdir(path):
52          os.makedirs(path)
53
54      # Node with task to extract all satelite bands
55      data_task = SentinelHubInputTask(
56          bands=bands,
57          bands_feature=(FeatureType.DATA, "BANDS"),
58          bands_dtype=float32,
59          data_collection=self.collection,
60          resolution=res,
61          maxcc=maxcc,
62          time_difference=datetime.timedelta(minutes=120),
63          config=self.config
64      )
65      data_node = EONode(data_task, name='data_node')
66
67      # Node with task to add nutrient values
68      groundtruth_values_task = AddFeatureTask((FeatureType.META_INFO, '
            GROUNDTRUTH'))
69      groundtruth_values_node = EONode(groundtruth_values_task, [data_node
            ],
70                                       name="groundtruth_values_node")
71      # Node with task to add date
72      date_task = AddFeatureTask((FeatureType.META_INFO, 'DATE'))
73      date_node = EONode(date_task, [groundtruth_values_node], name="
            date_node")
74
```

```python
75          # Node with task to save band names
76          bands_names_task = AddFeatureTask((FeatureType.META_INFO, 'BAND_NAMES
                '))
77          bands_names_node = EONode(bands_names_task, [date_node], name='
                bands_names_node')
78
79          # Node with task that saves the EOPatches != None into path
80          save = SaveNotNullTask(path, save_timestamps=True,
81                                 overwrite_permission=OverwritePermission.
                                        OVERWRITE_FEATURES)
82          save_node = EONode(save, [bands_names_node], name="save_node")
83
84          # Creates sequence of EONodes to be executed
85          workflow = EOWorkflow([data_node, groundtruth_values_node, date_node,
86                                 bands_names_node, save_node])
87
88          execution_args = []
89          for index, row in subset.iterrows():
90              # Gets the id
91              id = self.id if self.id is not None else index
92              # Calculates time interval
93              time_interval = get_time_interval(row['date'], padding_days,
                    date_fmt=date_fmt)
94
95              bbox = BBox(row['geometry'].bounds, crs=self.dataset.crs)
96              # Creates parameters for each EONode execution
97              execution_args.append({
98                  data_node: {'bbox': bbox, 'time_interval': time_interval},
99                  groundtruth_values_node: {'data': row[self.keep].to_dict()},
100                 date_node: {'data': row['date']},
101                 bands_names_node: {'data': self.bands},
102                 save_node: {'eopatch_folder': f'eopatch_{row[id]}'}
103             })
104
105         # Executes workflow
106         executor = EOExecutor(workflow, execution_args, save_logs=report)
107         executor.run()
108         # Saves report
```

```
109        if report:
110            executor.make_report()
```

**Listing B.32:** Code for tstk: SaveNotNullTask() in utils.eotasks.py

```python
1  class SaveNotNullTask(SaveTask):
2      """
3      #### Only saves EOPatches that have band data.
4      """
5      def execute(self, eopatch, **kwargs):
6          if eopatch[(FeatureType.DATA, 'BANDS')].size > 0:
7              super().execute(eopatch, **kwargs)
```

**Listing B.33:** Code for tstk: TSPatch() constructor

```python
1  class TSPatch(EOPatch):
2      def __init__(self, id: str, eopatch: EOPatch = None):
3          """
4          #### Extends the functionality of the original eo-patch
5              implementation with new methods
6
7          Args:
8              id: String used to identify the TSPatch
9              EOPatch: eo-learn abstraction to represent a single region
10          """
11         if eopatch is not None:
12             self.eopatch = eopatch
13         else:
14             super().__init__()
15             self.eopatch = super()
16         # Path for TSPatch in sys
17         self.id = id
18         # Saves location of EOPatch
19         self.location = self.bbox.geometry.centroid.x, self.bbox.geometry.
               centroid.y
20         # Saves available bands
21         self.bands, self.indices = {}, []
```

```
21          for v, b in enumerate(eopatch.meta_info['BAND_NAMES']):
22              self.bands[b] = v
```

**Listing B.34:** Code for tstk: TSPatch.add_indice()

```
1  def add_indice(self, indice: str, arr: np.ndarray):
2      """
3      #### Adds indice to eopatch
4
5      Args:
6          indice: name of indice to add
7          arr: np array with value to store
8      """
9      if type(indice) is not str:
10         raise TypeError(f'{indice} must be str. Got {type(indice)}')
11     if not isinstance(arr, np.ndarray):
12         raise TypeError(f'{arr} must be np.ndarray. Got {type(arr)}')
13     if np.ndim(arr) != 4:
14         raise ValueError(f'{arr} must have 4 dimentions (t, m, n, i). Got
               {arr.shape}')
15     if self.eopatch.__contains__((FeatureType.DATA, indice)):
16         raise ValueError(f'Indice {indice} already exists.')
17
18     self.data[indice] = arr
19     self.indices.append(indice)
```

**Listing B.35:** Code for tstk: TSPatch.get_location()

```
1  def get_location(self) -> Tuple[float, float]:
2      """
3      #### Getter for gps location of TSPatch
4
5      Returns:
6          Tuple: With latitude and longitude values
7      """
8      return self.location
```

**Listing B.36:** Prior code for tstk: TSPatch.get_masked_region()

```python
def get_masked_region(self):
    mask = self.patch.mask_timeless["IS_VALID"].squeeze()
    mask_filtered = np.where(mask==5,0,mask)
    return mask_filtered
```

**Listing B.37:** Code for tstk: TSPatch.get_masked_region()

```python
def get_masked_region(self, indice: str = None) -> np.ndarray:
    """
    #### Function that returns a numpy.array marking where data is valid
    #### When EOPatch has no mask info, all points are valid
    NOTE: 1 means valid. 0 means not valid


    Args:
        indice: String of indice to return array with it's shape. Default
            returns the bands mask


    Returns:
        numpy.ndarray: Mask to apply to patches, where 1 is valid
    """
    # All points are valid if there is no mask
    if not self.eopatch.__contains__(FeatureType.MASK_TIMELESS):
        if indice is not None:
            # Checks for indice existance
            if not self.eopatch.__contains__((FeatureType.DATA, indice)):
                raise ValueError(f'Indice {indice} does not exist. Only {
                    self.eopatch.get_features()}')


            # Returns with indice's shape full of 1's
            return np.ones_like(self.data[indice][0])


        # Returns array like BANDS full of 1's
        return np.ones_like(self.data['BANDS'][0])


    mask = self.mask_timeless["IS_VALID"].squeeze()
    mask_filtered = np.where(mask==5, 0, mask)
```

```
28          return mask_filtered
```

**Listing B.38:** Prior code for tstk: TSPatch.get_masked_region_values()

```python
1  def get_values_of_masked_region(self ,indices = None ,band_names = None,
       as_array=True):
2      """Returns the pixels in the masked region for the selected `indices`
            and `band_names` for each of the patch
3
4      Args:
5          indices (str): The list of indices in which we want to get the
                values of
6          band_names (str): The list of indices in which we want to get the
                 values of
7          as_array (bool, optional): If true returns in 1D array form(only
                the values with data), else returns in 2D array. Defaults to
                True.
8
9      Returns:
10         ndarray: If `as_array` is true returns in 1D array form(only the
                values with data), else returns in 2D array.
11     """
12     values = {}
13     eopatch = self.patch
14     masked_region = self.get_masked_region()
15     nearest_image_index = self._get_index_nearest_to_collection_date()
16     for i in range(0,eopatch.data["BANDS"].shape[-1]):
17         values[band_names[i]] = eopatch.data["BANDS"][nearest_image_index
                ,...,i]*masked_region
18         if as_array:
19             values[band_names[i]] = values[band_names[i]][masked_region!=
                    0]
20
21     for i in indices:
22         values[i] = eopatch.data[i][nearest_image_index ,...,-1]*
                masked_region
23         if as_array:
24             values[i] = values[i][masked_region!=0]
```

```
25          return values
```

**Listing B.39:** Code for tstk: TSPatch.get_masked_region_values()

```python
1  def get_values_of_masked_region(self, indices: List[str] = None, bands: List[
       str] = None,
2                                   as_array: bool = True, verbose: bool =
                                        False
3                                   ) -> dict[str: np.ndarray]:
4       """
5       #### Returns the pixels in the masked region for the selected `
           indices` and `bands`\
6       for each of the patch
7
8       Args:
9           indices: The list of indices in which we want to get the values
               of.
10                   None means all indices.
11          bands: The list of bands names of which we want to get the values
                of.
12                  None means all bands.
13          as_array: If true returns in 1D array form(only the values with
               data),
14                   else returns in 2D array. Defaults to True.
15          verbose: boolean to print important runtime information.
16
17      Returns:
18          dict: Dictionary mapping indice/band name to numpy.ndarray with
               values
19      """
20      result = {}
21
22      # Gets masked region for indices
23      if indices is not None:
24          # Indices must be list
25          if type(indices) is not list:
26              raise TypeError(f'Indices must be list. Got {type(indices)}')
27          # All in indices must be str
```

```python
28              if not all(type(val) is str for val in indices):
29                  raise TypeError(f'Indices must be list of strings. Got {type(
                        indices[0])}')
30              # Check if indices exist in eopatch
31              if not set(indices).issubset(self.indices):
32                  raise ValueError(f'Indice must exist. Available are: {self.
                        indices}')

33

34          indices = self.indices if indices is None else indices

35

36          # Parser input band names
37          if bands is not None:
38              # Bands must be list
39              if type(bands) is not list:
40                  raise TypeError(f'Bands must be list. Got {type(bands)}')
41              # All in bands must be str
42              strings = [type(b) is str for b in bands]
43              if not all(strings):
44                  wrong_idx = strings.index(False)
45                  raise TypeError(f'Bands must be list of strings. Got {type(
                        bands[wrong_idx])}')
46              # Check if bands exist in eopatch
47              if not set(bands).issubset(self.bands.keys()):
48                  raise ValueError(f'Band must exist. Available are: {self.
                        bands.keys()}')

49

50          bands = list(self.bands.keys()) if bands is None else bands

51

52          # Closest to date index
53          closest_i = self._get_index_nearest_to_collection_date(verbose)

54

55          if len(indices) > 0:
56              # Gets mask with indice shapelike
57              indices_mask = self.get_masked_region(indices[0])
58              # Calculate mask for indices
59              for indice in indices:
60                  result[indice] = self.data[indice][closest_i] * indices_mask
61                  if as_array:
```

```
62                      result[indice] = result[indice].flatten()

63

64           # Gets masked region for bands
65           for band in bands:
66               idx = self.bands[band]
67               # Gets mask only for current band
68               band_mask = self.get_masked_region()[..., idx]
69               result[band] = self.data['BANDS'][closest_i][..., idx] *
                     band_mask
70               if as_array:
71                   result[band] = result[band].flatten()

72

73           return result
```

**Listing B.40:** Prior code for tstk: TSPatch.represent_image()

```python
1  def represent_image(self,estimation):
2          """Draws an image with the values estimated
3
4          Args:
5              estimation (array): Array with the size of the masked region(1D)
6
7          Returns:
8              array: 2D image
9          """
10         nearest_image_index = self._get_index_nearest_to_collection_date()
11
12         eopatch = self.patch
13         mask = self.get_masked_region()
14         mask = mask.astype(float)
15         image = eopatch.data["BANDS"][nearest_image_index][...,[3,2,1]]
16         _max = 255#dfeopatches["N"].max()
17         _min = 0#dfeopatches["N"].min()
18         convert_est = estimation/_max
19         for i,tup in enumerate(zip(mask.nonzero()[0],mask.nonzero()[1])):
20             x = tup[0]
21             y = tup[1]
22             try:
```

```
23                          image[x,y] = [convert_est[i], 0,0]
24                  except BaseException as err:
25                      print(f"The estimation most likely doesn't correspond to the
                             eopatch: {err}")
26          return image
```

**Listing B.41:** Code for tstk: TSPatch.represent_image()

```python
1  def represent_image(self, R: str, G: str, B: str, factor: int = 255, verbose:
       bool = False
2                         ) -> np.ndarray:
3          """
4          #### Returns an image array in RGB with values from features
               mentioned.
5
6          Args:
7              R: Name of indice/band used for Red axis.
8              G: Name of indice/band used for Green axis.
9              B: Name of indice/band used for Blue axis.
10             factor: Float to multiply with values
11             verbose: boolean to print important runtime information.
12
13         Returns:
14             array: 2D image
15         """
16         # Gets the index of the patch taken closer to the rewuested time
17         nearest_indx = self._get_index_nearest_to_collection_date(verbose)
18
19         # Saves bands shape
20         bands_shape = self.data['BANDS'].shape
21         # Image is line, column and [R, G, B]
22         img = np.ones(shape=(bands_shape[1], bands_shape[2], 3))
23
24         # Check arguments
25         for n, axis in enumerate((R, G, B)):
26             # Axis is a prev calculated indice
27             if type(axis) is str:
28                 if axis not in self.indices and axis not in self.bands.keys()
```

```
              :
29               raise ValueError(f'Parameter does not exist. Available: {
                     self.indices} or {list(self.bands.keys())}')

30

31           if axis in self.bands.keys():
32               # Get Band values
33               img[..., n] = self.data['BANDS'][nearest_indx][..., self.
                     bands[axis]] * factor
34           else:
35               # Get indice values
36               img[..., n] = self.data[axis][nearest_indx][..., 0] *
                     factor

37

38       else:
39           raise TypeError(f'Parameters must be str. Got: {type(R), type
                 (G), type(B)}')

40

41       return img
```

**Listing B.42:** Prior code for tstk: TSPatch._get_index_nearest_to_collection_date()

```
1   def _get_index_nearest_to_collection_date(self):
2       smallest_index = 0
3       smallest_difference = dt.timedelta(days=2000)
4       try:
5           collected_day = dt.datetime.strptime(self.get_dataset_entry_value
                 ("SURVEY_DATE"),'%d/%m/%y')
6       except:
7           return -1
8       for i,image_date in enumerate(self.timestamp):
9           current_difference = abs(collected_day - image_date)
10          if(current_difference < smallest_difference):
11              smallest_difference = current_difference
12              smallest_index = i
13      return smallest_index
```

**Listing B.43:** Code for tstk: TSPatch._get_index_nearest_to_collection_date()

```python
def _get_index_nearest_to_collection_date(self, verbose: bool = False):
    """
    This function gets the index of the temporal band closest to the
        collection data that is not\
    full of 0's.


    Args:
        verbose: boolean to print important runtime information.

    Raises ValueError if all temporal bands are full of 0's
    """
    collected_day = datetime.strptime(self.meta_info['date'], '%d/%m/%Y')

    if len(self.timestamp) == 1:
        result = 0
    else:
        sorted_timestamps = sorted(self.timestamp, key=lambda t: abs(
            collected_day - t))
        result = self.timestamp.index(sorted_timestamps[0])

    # Check if array is full of zeros across all dimensions
    if np.all(self.data['BANDS'][result] == 0, axis=None):
        if verbose:
            print('INFO: EOP nearest to collection date is full of 0.')
        if len(self.timestamp) == 1:
            raise ValueError(f'INFO: No more EOP in TSP, thus ignoring
                this {self.id}.')

        if verbose:
            print('INFO: Getting 2nd closest and non-zero EOP.')
        for t in range(1, len(sorted_timestamps)):
            indx = self.timestamp.index(sorted_timestamps[t])
            if np.all(self.data['BANDS'][indx] != 0, axis=None):
                if verbose:
                    print(f'INFO: Got EOP with {abs(collected_day -
                        sorted_timestamps[t])} days different.')
                return indx
```

```
34
35          raise ValueError(f'INFO: No more EOP in TSP, thus ignoring this {
                self.id}.')

36

37      return result
```

**Listing B.44:** Code for tstk: Dataset() and Dataset.__create_append_tsp__()

```python
1  class Dataset:
2      """
3      #### Class that represents a collection of EOPatches saved in disk
4
5      Args:
6          eops_path: String with path to directory with EOPatches to be read
7      """
8      def __init__(self, eops_path: str):
9          # Checks if path exists
10          if type(eops_path) is not str:
11              raise TypeError(f'Path must be str.')
12          if not isdir(eops_path):
13              raise ValueError(f'Path given: {eops_path} is not directory.')
14
15          # Path to directory with eopatches
16          self.eops_path = eops_path
17
18          # Parallelize TSP creation
19          self.tspatches = []
20          self.lock = Lock()
21          list_dir = listdir(eops_path)
22          with concurrent.futures.ThreadPoolExecutor() as executor:
23              executor.map(self.__create_append_tsp__, list_dir)
24
25          # Pandas dataframe with info about TSPatches groundtruth data
26          self.gth_df = None
27          # Saves available indices in EOPatches
28          self.indices = []
29          # Saves dict of available bands in EOPatches
30          self.bands = self.tspatches[0].get_bands()
```

```
31          # Saves groundtruth data names
32          self.groundtruth_f = self.tspatches[0].get_groundtruth_features()
33
34
35      # Method to parallelize TSP creation
36      def __create_append_tsp__(self, eop_dir):
37          full_path = join(self.eops_path, eop_dir)
38          tsp = TSPatch.load(full_path, lazy_loading=True)
39          with self.lock:
40              self.tspatches.append(tsp)
```

**Listing B.45:** Code for tstk: Dataset.add_indices()

```
1 def add_indices(self, indice_func: dict[str: Callable], save: bool = False,
    report: bool = False,
2                    verbose: bool = False, override_values: bool = False):
3          """
4          #### Function that receives a mapping for indice's names to how they
               are calculated, \
5          and adds them to list of TSPatches
6
7          Args:
8              indice_func: Dictionary mapping indice name and how it is
                   calculated.
9              save: Boolean to overwrite data and save indice to EOPatch in sys
                   .
10             report: Boolean specifies if user wants execution report.
11             verbos: Boolean to show important runtime information.
12             override_values: When True, tstk changes the band values that
                   give mathematical errors.
13
14         Example:
15             >>> bands = dataset.get_bands()
16             >>> indices = {'EVI2': lambda b: 2.4*(b[bands['B:NIR']] - b[bands
                   ['B:GREEN']]}
17             >>> dataset.add_indices(indices)
18             >>> dataset.add_indices(indices, save=True, report=True)
19
```

**167**

```python
20          NOTE: All temporal bands full of 0 are ignored. Moreover, if valid
                temporal band\
21              has 0 values, they are incremented by 1e-3 to avoid mat errors.
22          """
23          np.seterr(all='raise')

24

25          # Node to save new indice in
26          if save:
27              save_task = SaveTask(self.eops_path, features=[(FeatureType.DATA,
                    indice)])
28              save_node = EONode(save_task, name=f'save_{indice}_node')
29              workflow, exec_args = EOWorkflow([save_node]), []

30

31          for tsp in self.get_tspatches():
32              eop_bands = tsp.data['BANDS']

33

34              # invalid temporal bands <=> bands with nothing but 0
35              inv_t_b = np.all(eop_bands==0, axis=(1, 2, 3))
36              # Some or all temporal bands are full of 0
37              if np.sum(inv_t_b) > 0:
38                  if verbose:
39                      print(f'INFO: TSP {tsp.id} has {np.sum(inv_t_b)} temporal
                            bands full of 0. These are ignored and kept as 0.')
40                  # All temporal bands are invalid. Skipping to next tsp
41                  if np.sum(inv_t_b) == eop_bands.shape[0]:
42                      if verbose:
43                          print(f'INFO: Skipping TSP {tsp.id} calculation. All
                                temporal bands are 0.')
44                      continue

45

46              # Increments random small value on all valid temporal bands to
                    avoid calculation errors
47              if np.sum(eop_bands[~inv_t_b] == 0) > 0:
48                  if verbose:
49                      print(f'INFO: TSP {tsp.id} has {np.sum(eop_bands[~inv_t_b
                            ] == 0)} pixels = 0 in valid temporal bands. These
                            are incremented by 1e-3.')
50                  zero_indices = np.where(eop_bands[~inv_t_b] == 0)
```

```python
51                 eop_bands[zero_indices] += np.random.random(size=len(
                        zero_indices[0])) * 1e-3
52
53            for indice, func in indice_func.items():
54
55                if tsp.__contains__((FeatureType.DATA, indice)):
56                    if verbose:
57                        print(f'INFO: TSP {tsp.id} already has indice {indice
                            }. Skipping calculation.')
58                    continue
59
60                indice_eop = np.zeros(shape=eop_bands.shape[:-1], dtype=np.
                        float32)
61                # Apply function. func is applied to each 3rd dimetion = Each
                        pixel
62                try:
63                    indice_eop[~inv_t_b] = np.apply_along_axis(func, axis=3,
                            arr=eop_bands[~inv_t_b])
64                except FloatingPointError:
65                    if verbose:
66                        print(f'SVI {indice} is mathematically impossible to
                            calculate on TSP {tsp.id}')
67
68                # Bands temporal is (n, m) but to store as DATA it needs to
                        be (n, m, b)
69                # even though it should be 3d, storing as DATA_TIMELESS seems
                        incorrect
70                # So we add an unecessary dimention
71                indice_eop = np.expand_dims(indice_eop, axis=-1)
72                tsp.add_indice(indice, indice_eop)
73
74            if save:
75                exec_args.append({save_node: {'eopatch_folder': tsp.id, '
                        eopatch': tsp.eopatch}})
76
77        # Save in sys
78        if save:
79            executor = EOExecutor(workflow, exec_args, save_logs=report)
```

```
80          executor.run(multiprocess=False)
81          # Saves report
82          if report:
83              executor.make_report()
84
85      for i in indice_func.keys():
86          if i not in self.indices:
87              self.indices.append(i)
```

**Listing B.46:** Prior code for tstk: Dataset.add_index() and Dataset.save_indices_to_eopatches()

```
1 def add_index(self, index_name, index_formula):
2       """Adds a specific index to the dataset in question
3
4       NOTE:To make the index persist `save_indices_to_patches` must be
            called
5
6       Args:
7           index_name (str): The name of the index
8           index_formula (str): The formula for the specified index
9
10      Example:
11
12          >>> dataset.add_index("NDVI","(B07-B04)/(B07+B04)")
13          >>> dataset.add_index("IRECI","(B07-B04)/(B05/B06)")
14          >>> dataset.save_indices_to_patches()
15
16      """
17      self.index_dic[index_name] = index_formula
18
19
20
21  def save_indices_to_eopatches(self):
22      """Saves the previously indicated indices into the dataset.
23
24      Example:
25
26          >>> dataset.add_index("NDVI","(B07-B04)/(B07+B04)")
```

```
27              >>> dataset.add_index("IRECI","(B07-B04)/(B05/B06)")
28              >>> dataset.save_indices_to_patches()
29          """
30          load = LoadTask(self.eopatches_folder)
31          load_node = EONode(load,name="Load")
32          available_bands = ['B01', 'B02', 'B03', 'B04', 'B05', 'B06', 'B07', '
                B08', 'B8A', 'B09', 'B10', 'B11', 'B12']
33
34          add_indices = AddIndicesTask(self.index_dic,available_bands)
35          add_indices_node = EONode(add_indices,[load_node],name="Add_index")
36          save = SaveTask(self.eopatches_folder, overwrite_permission=
                OverwritePermission.OVERWRITE_PATCH)
37          add_save_node = EONode(save,[add_indices_node],name="Save_patch")
38          execution_args = []
39          workflow = EOWorkflow([load_node, add_indices_node,add_save_node])
40          eopatch_folders = os.listdir(self.eopatches_folder)
41          print(len(eopatch_folders))
42          for i in eopatch_folders:
43              execution_args.append(
44              {
45                  load: {'eopatch_folder': f'{i}'},
46                  add_indices: {},
47                  save: {'eopatch_folder': f'{i}'}
48              })
49          executor = EOExecutor(workflow, execution_args, save_logs=True)
50          executor.run(workers=5, multiprocess=False)
51          executor.make_report()
```

**Listing B.47:** Code for tstk: Dataset.add_indices()

```
1 def add_indices(self, indice_func: dict[str: Callable], save: bool = False,
     report: bool = False,
2                  verbose: bool = False, override_values: bool = False):
3          """
4          #### Function that receives a mapping for indice's names to how they
                are calculated, \
5          and adds them to list of TSPatches
6
```

```python
 7          Args:
 8              indice_func: Dictionary mapping indice name and how it is
                     calculated.
 9              save: Boolean to overwrite data and save indice to EOPatch in sys
                     .
10              report: Boolean specifies if user wants execution report.
11              verbos: Boolean to show important runtime information.
12              override_values: When True, tstk changes the band values that
                     give mathematical errors.
13
14          Example:
15              >>> bands = dataset.get_bands()
16              >>> indices = {'EVI2': lambda b: 2.4*(b[bands['B:NIR']] - b[bands
                     ['B:GREEN']]}
17              >>> dataset.add_indices(indices)
18              >>> dataset.add_indices(indices, save=True, report=True)
19
20          NOTE: All temporal bands full of 0 are ignored. Moreover, if valid
                 temporal band\
21              has 0 values, they are incremented by 1e-3 to avoid mat errors.
22          """
23          np.seterr(all='raise')
24
25          # Node to save new indice in
26          if save:
27              save_task = SaveTask(self.eops_path, features=[(FeatureType.DATA,
                     indice)])
28              save_node = EONode(save_task, name=f'save_{indice}_node')
29              workflow, exec_args = EOWorkflow([save_node]), []
30
31          for tsp in self.get_tspatches():
32              eop_bands = tsp.data['BANDS']
33
34              # invalid temporal bands <=> bands with nothing but 0
35              inv_t_b = np.all(eop_bands==0, axis=(1, 2, 3))
36              # Some or all temporal bands are full of 0
37              if np.sum(inv_t_b) > 0:
38                  if verbose:
```

```python
39                    print(f'INFO: TSP {tsp.id} has {np.sum(inv_t_b)} temporal
                          bands full of 0. These are ignored and kept as 0.')
40              # All temporal bands are invalid. Skipping to next tsp
41              if np.sum(inv_t_b) == eop_bands.shape[0]:
42                  if verbose:
43                      print(f'INFO: Skipping TSP {tsp.id} calculation. All
                              temporal bands are 0.')
44                  continue
45
46          # Increments random small value on all valid temporal bands to
                avoid calculation errors
47          if np.sum(eop_bands[~inv_t_b] == 0) > 0:
48              if verbose:
49                  print(f'INFO: TSP {tsp.id} has {np.sum(eop_bands[~inv_t_b
                          ] == 0)} pixels = 0 in valid temporal bands. These
                          are incremented by 1e-3.')
50              zero_indices = np.where(eop_bands[~inv_t_b] == 0)
51              eop_bands[zero_indices] += np.random.random(size=len(
                    zero_indices[0])) * 1e-3
52
53          for indice, func in indice_func.items():
54
55              if tsp.__contains__((FeatureType.DATA, indice)):
56                  if verbose:
57                      print(f'INFO: TSP {tsp.id} already has indice {indice
                              }. Skipping calculation.')
58                  continue
59
60              indice_eop = np.zeros(shape=eop_bands.shape[:-1], dtype=np.
                    float32)
61              # Apply function. func is applied to each 3rd dimetion = Each
                     pixel
62              try:
63                  indice_eop[~inv_t_b] = np.apply_along_axis(func, axis=3,
                        arr=eop_bands[~inv_t_b])
64              except FloatingPointError:
65                  if verbose:
66                      print(f'SVI {indice} is mathematically impossible to
```

```python
                            calculate on TSP {tsp.id}')

                    # Bands temporal is (n, m) but to store as DATA it needs to
                        be (n, m, b)
                    # even though it should be 3d, storing as DATA_TIMELESS seems
                        incorrect
                    # So we add an unecessary dimention
                    indice_eop = np.expand_dims(indice_eop, axis=-1)
                    tsp.add_indice(indice, indice_eop)


            if save:
                exec_args.append({save_node: {'eopatch_folder': tsp.id, '
                    eopatch': tsp.eopatch}})


        # Save in sys
        if save:
            executor = EOExecutor(workflow, exec_args, save_logs=report)
            executor.run(multiprocess=False)
            # Saves report
            if report:
                executor.make_report()


        for i in indice_func.keys():
            if i not in self.indices:
                self.indices.append(i)
```