

```

1 ######
2 # Created by Ben Marosites
3 # marosite@email.sc.edu
4 # For GEOL 365/599
5 # Fall 2025 - December 01
6 # These are functions I created or modified to use in my analysis.
7 # Most is related to extracting and processing data.
8 #####
9
10 import pandas as pd
11 import numpy as np
12 import math
13 import datetime as dt
14 from matplotlib import pyplot as plt
15
16 import os
17 import shutil
18 import glob
19
20 from netCDF4 import Dataset
21 import xarray as xr
22
23 from scipy.interpolate import griddata
24 from scipy.signal import argrelextrema, find_peaks
25
26 def extract_tar_files(directory):
27     """
28         This was just cut from one of my notebooks. Should only be run once to extract data
29         and move it to the correct final location. This function will need to be redone, but
30         for now, it will be incomplete.
31
32     Parameters
33     -----
34
35     Returns
36     -----
37
38     Examples
39     -----
40
41     """
42     extract = False # ONLY NEED TO RUN THIS ONCE, SET TO FALSE AFTER DATA IS IN THE
43     # CORRECT DIR.
44     if extract:
45         os.makedirs('./data/raw_lidar/extracted_data', exist_ok=True)
46         # dir = './data/raw_lidar'
47         os.listdir(directory)
48         extracted_lidar_dir = './data/raw_lidar/extracted_data'
49         os.makedirs('./data/raw_lidar/extracted_data', exist_ok=True)
50
51         for file in os.listdir(directory):
52             # we have tar files. we can extract these with a lf function
53             if file.endswith('.tar'):
54                 full_path_file = os.path.join(directory, file)
55                 print(f'extracting {full_path_file} to {extracted_lidar_dir}...')
56                 lf.extract_data(full_path_file, extracted_lidar_dir)
57
58             # Tar files are files of files. This creates a messy directory. This will move
59             # it from that directory to make it a little easier to navigate
60             # Define the source and destination paths
61             temp_dir = './data/raw_lidar/extracted_data/Users/HISCOX/OneDrive - University
62             of South Carolina/Research/SAVANT/SAVANT Data/Level 1 - Grouped/USC Aerosol
63             Lidar/ASCII'
64             source_directory = './data/raw_lidar/extracted_data/ASCII/'
65             destination_directory = extracted_lidar_dir
66
67             try:

```

```

4 # Move the directory
5     for folder in os.listdir(source_directory):
6         print(folder)
7         shutil.move(os.path.join(source_directory, folder), destination_directory)
8         print(f"Directory '{source_directory}' moved successfully to '{destination_directory}'")
9 except FileNotFoundError:
10     print(f"Error: Source directory '{source_directory}' not found.")
11 except Exception as e:
12     print(f"An error occurred: {e}")
13
14
15 def read_fastdata(file):
16     """
17     Read high-frequency data from a netCDF file into a pandas dataframe.
18     This is developed specifically for the SAVANT wind data I have available.
19     This data is collected at 4-hour increments, with a sampling rate of 20 Hz (20 times
20     per second).
21     I am interested in wind and temperature data to decompose the heat flux.
22
23     Parameters
24     -----
25
26     Returns
27     -----
28
29     Examples
30     -----
31
32     """
33
34     with Dataset(file, 'r') as ds:
35         time_units = ds.variables['time'].units[14:]
36         time_1s = ds.variables['time'][:]
37         num_intervals = ds.dimensions['time'].size    # should be 14400, for 4 hours in
38         seconds
39         samples_per_interval = ds.dimensions['sample'].size    # should be 20 for 20
40         hertz or 20 times per second
41         total_samples = num_intervals * samples_per_interval
42
43         delta_t = ds.variables['time'].__dict__['interval(sec)']/samples_per_interval
44
45         start_time = time_1s[0]
46         end_time = start_time + 14400
47
48         time_20hz = np.arange(start_time, end_time, delta_t)
49         print(len(time_20hz))
50
51         # extract wind data
52         data_dict = {}
53         wind_fields = [var for var in ds.variables.keys() if hasattr(ds.variables[var], 'long_name')]
54             and 'wind' in ds.variables[var].long_name.lower()]
55         for field in wind_fields:
56             if ds.variables[field].shape == (14400, 20):
57                 samples_2d = ds.variables[field]::
58                 time_series_1d = samples_2d.flatten().data
59                 data_dict[field] = time_series_1d
60             else:
61                 print(f'{field} not used. Incorrect length.')
62         temperature_fields = [var for var in ds.variables.keys() if hasattr(ds.variables[var], 'long_name')]
63             and 'temp' in ds.variables[var].long_name.lower()]
64         for field in temperature_fields:
65             if ds.variables[field].shape == (14400, 20):
66                 samples_2d = ds.variables[field]::
67                 time_series_1d = samples_2d.flatten().data

```

```

125         data_dict[field] = time_series_1d
126     else:
127         print(f'{field} not used. Incorrect length.')
128
129     df = pd.DataFrame(data_dict, index=time_20hz)
130
131     try:
132         df.index = pd.to_datetime(df.index, unit='s', origin=time_units)
133     except ValueError as e:
134         print(f"Warning: Could not set datetime index due to error: {e}")
135         # If conversion fails, keep the index as seconds (float)
136
137     df.index.name = 'time'
138
139     return df
140
141 def calc_wind(u, v):
142     """
143         Calculates wind speed (m/s) and direction (0-360 degrees) given
144         u (zonal) and v (meridional) components.
145
146         Parameters
147         -----
148             u (float, int): wind speed in zonal direction or x
149             v (float, int): wind speed in meridional direction or y
150
151         Returns
152         -----
153             windspeed and direction
154
155         Examples
156         -----
157             spd, dir = calc_wind(3,4)
158             print(spd, dir)
159
160
161     """
162     speed = np.sqrt(u**2 + v**2)
163     direction = (np.rad2deg(np.arctan2(u, v)) + 360) % 360
164
165     return speed, direction
166
167 def create_winds(df):
168     """
169         Adds wind speed and wind direction using the calc_wind function.
170
171         Parameters
172         -----
173
174         Returns
175         -----
176
177         Examples
178         -----
179
180     """
181     df = df.copy()
182     height_suffixes = set()
183     for col in df.columns:
184         suffix = col[col.find('_'): ]
185         print(col, suffix)
186         height_suffixes.add(suffix)
187
188     for suffix in height_suffixes:
189         u_col = 'u'+suffix
190         v_col = 'v'+suffix
191         if u_col in df.columns and v_col in df.columns:
192             print(f'calculating {u_col} and {v_col}')

```

```

192     df[f'spd{suffix}'], df[f'dir{suffix}'] = calc_wind(df[u_col], df[v_col])
193
194     return df
195
196
197     def read_metadata(directory, file_path = None):
198         """
199             Reads metadata from Raymetrics lidar .txt files.
200
201             Parameters
202             -----
203                 directory (string): directory of lidar txt files.
204                 file_path (string, optional): path to lidar file. if None, the entire directory
205                 will be read.
206
207             Returns
208             -----
209                 Pandas DataFrame with metadata from txt file(s) processed.
210
211             Examples
212             -----
213
214         if file_path == None:
215             files_to_read = [txt for txt in os.listdir(directory) if txt.endswith('.txt')]
216         if file_path != None:
217             files_to_read = [file_path]
218
219         meta_df = pd.DataFrame(columns=['file', 'start', 'end', 'elevation', 'longitude',
220                                       'latitude', 'zenith', 'azimuth', 'temp_ground', 'pressure'])
221
222         # Convert appropriate columns to numeric
223         numeric_cols = ['elevation', 'longitude', 'latitude', 'zenith', 'azimuth',
224                         'temp_ground', 'pressure']
225         # meta_df[numeric_cols] = meta_df[numeric_cols].apply(pd.to_numeric,
226         # errors='coerce') # Convert and set non-numeric values to NaN
227
228         for txt_file in files_to_read:
229             f = os.path.join(directory, txt_file)
230             meta_data = []
231             shot_data = []
232             with open(f, 'r') as file:
233                 content = file.readlines()
234
235                 # The first seven lines are meta data
236                 meta = content[0:7]
237
238                 header1 = meta[0]
239                 header2 = meta[1]
240                 header3 = meta[2]
241                 header4 = meta[3]
242                 header5 = meta[4]
243                 header6 = meta[5]
244                 header7 = meta[6]
245
246                 header2 = header2.split(' ')
247                 campaign = header2[0]
248                 startDate = f'{header2[3]} {header2[4]}'
249                 endDate = f'{header2[5]} {header2[6]}'
250                 elevation = float(header2[7])
251                 longitude = float(header2[8])
252                 latitude = float(header2[9])
253                 zenith = float(header2[10])*-1
254                 azimuth = float(header2[11])
255                 temp_ground = float(header2[12])
256                 pressure = float(header2[13])
257                 bin_width = float(header5.split()[6])
258
259                 data = content[7:]

```

```

255
256     # Create a dataframe for the meta data
257     meta_data.append([f, startDate, endDate, elevation, longitude, latitude, zenith,
258                       azimuth, temp_ground, pressure])
259
260     # Create DF
261     temp_meta_df = pd.DataFrame(meta_data, columns=['file', 'start', 'end',
262                                   'elevation', 'longitude', 'latitude', 'zenith', 'azimuth', 'temp_ground',
263                                   'pressure'])
264
265     # Convert appropriate columns to numeric
266     # numeric_cols = ['elevation', 'longitude', 'latitude', 'zenith', 'azimuth',
267     #                  'temp_ground', 'pressure']
268     # temp_meta_df[numeric_cols] = temp_meta_df[numeric_cols].apply(pd.to_numeric,
269     # errors='coerce') # Convert and set non-numeric values to NaN
270
271     # Convert date columns to datetime
272     temp_meta_df['start'] = pd.to_datetime(temp_meta_df['start'], errors='coerce',
273                                         dayfirst=True)
274     temp_meta_df['end'] = pd.to_datetime(temp_meta_df['end'], errors='coerce',
275                                         dayfirst=True)
276     meta_df = pd.concat([meta_df, temp_meta_df], ignore_index=True)
277     meta_df = meta_df.sort_values('start')
278
279     return meta_df
280
281
282     def find_scans(df, start_time=None, end_time=None):
283         """
284         """
285         # find local min max indices
286         max_zenith_idxes= find_peaks(df['zenith'].values)[0]
287         #argrelextrema(df['zenith'].values, np.greater)[0]
288         min_zenith_idxes = find_peaks(-df['zenith'].values)[0] #
289         argrelextrema(df['zenith'].values, np.less)[0]
290
291         # Merge and sort extrema indices
292         min_max = np.sort(np.concatenate((max_zenith_idxes, min_zenith_idxes)))
293         min_max = np.insert(min_max, 0, 0) # insert 0 so we can start with the initial row.
294         if len(df['zenith'])-1 not in min_max: # add last row to the end.
295             min_max = np.insert(min_max, len(min_max), len(df['zenith'])-1)
296
297         # organize lists of indexes
298         times = []
299         zeniths = []
300         for extrema in min_max:
301             row = df.iloc[extrema]
302             times.append(row['start'])
303             zeniths.append(row['zenith'])
304
305             start_stop_pairs = []
306             for i in range(len(min_max)-1):
307                 if min_max[i+1] - min_max[i] > 4:
308                     start_stop_pairs.append([int(min_max[i]), int(min_max[i+1])])
309
310             time_stamp_pairs = []
311             for pair in start_stop_pairs:
312                 time_stamp_pairs.append([df.iloc[pair[0]]['start'], df.iloc[pair[1]]['start']])
313
314             # Plot a figure to test.
315             plt.figure(figsize=(18,8))
316             plt.plot(df['start'], df['zenith'])
317             plt.title(f"{df.iloc[0]['start']}-{df.iloc[-1]['start']}")
318
319             plt.scatter(df['start'], df['zenith'], s=10)
320             plt.scatter(times, zeniths, color='red')
321
322             if start_time != None and end_time != None:

```

```

313     t1 = pd.to_datetime(start_time)
314     t2 = pd.to_datetime(end_time)
315     plt.xlim(t1, t2)
316     for row in start_stop_pairs:
317         plt.axvspan(df.iloc[row[0]]['start'], df.iloc[row[1]]['start'], color='grey',
318                     alpha=0.25)
319     plt.show()
320     return start_stop_pairs, time_stamp_pairs
321
322 def file_to_data(txt_file, max_distance=400, bin_width=7.5, origin_X=0, origin_Y=0,
323                  origin_Z=2, origin_azimuth=0, origin_zenith=0, bg_len_cutoff=4, filter_t=2):
324     """
325         Reads a LiDAR text file, extracts measurement data, processes background noise
326         correction, and computes
327         spatial coordinates for each data point.
328
329     Parameters
330     -----
331         txt_file (str): Path to the LiDAR text file.
332         distance (float, optional): Maximum measurement distance. Default is 450.
333         bin_width (float, optional): Width of each measurement bin. Default is 7.5.
334         origin_X (float, optional):
335         origin_Y (float, optional):
336         origin_Z (float, optional):
337         origin_azimuth (float, optional):
338         origin_zenith (float, optional):
339         filter_t (float, optional): threshold filter. Default set to 1 * sigma
340
341     Returns
342     -----
343         pandas.DataFrame: A DataFrame containing the processed LiDAR data with background
344         noise correction and
345             spatial coordinates (x, z, distance).
346
347     Examples
348     -----
349
350     """
351     with open(txt_file, 'r') as file:
352         content = file.readlines()
353
354         meta = content[:7]
355         shot_data = meta[1]
356         Location, StartDate, StartTime, EndDate, EndTime, Elevation, Longitude, Latitude,
357         Zenith, Azimuth, Temp, Pressure = shot_data.split()
358         Zenith = -float(Zenith)
359         last_row = int(math.ceil(max_distance/bin_width))
360         data = [row.strip().split('\t') for row in content[7:(7+last_row)]]
361
362         data = pd.DataFrame(data, columns=['analog', 'photon']).apply(pd.to_numeric)
363
364         # Add data from meta data
365         data['start'] = pd.to_datetime(f'{StartDate} {StartTime}', format='%d/%m/%Y %H:%M:%S')
366         data['end'] = pd.to_datetime(f'{EndDate} {EndTime}', format='%d/%m/%Y %H:%M:%S')
367         data['zenith'] = Zenith
368
369         # Add distance information
370         step_size = bin_width
371         data['distance'] = (data.index+1) * step_size
372         data['x'] = data['distance'] * np.cos(np.radians(Zenith))
373         data['z'] = origin_Z + data['distance'] * np.sin(np.radians(Zenith))
374
375         # First perform Background correction. We need to remove the noise.
376         bg_length = min(1000, int(len(data)/bg_len_cutoff))    # bg_length = min(1000,
377                      int(len(data)/4))

```

```

373 EndSig_A, EndSig_P = data['analog'][-bg_length:].mean(), data['photon'][-bg_length:].
374     mean()
375 EndSig_A_std, EndSig_P_std = data['analog'][-bg_length:].std(), data['photon'][-bg_length:].
376     std()
377
378     data['analog_bgc'] = np.where(data['analog'] >= EndSig_A + filter_t * EndSig_A_std,
379         data['analog'] - EndSig_A, 0) # old:---> data['analog_bgc'] =
380         np.where(data['analog'] >= EndSig_A + 3 * EndSig_A_std, data['analog'] - EndSig_A, 0)
381     data['photon_bgc'] = np.where(data['photon'] >= EndSig_P + filter_t * EndSig_P_std,
382         data['photon'] - EndSig_P, 0) # old:---> data['photon_bgc'] =
383         np.where(data['photon'] >= EndSig_P + 3 * EndSig_P_std, data['photon'] - EndSig_P, 0)
384
385     # Range correction. Inverse-square law of light.
386     Ranges=data['distance'].values
387     data['analog_rcs']=data['analog_bgc'] *Ranges**2
388     data['photon_rcs']=data['photon_bgc'] *Ranges**2
389
390     # Normalize data
391     peak=max(data['analog_rcs'])
392     data['analog_rcs_norm']=data['analog_rcs']/peak
393
394     return data
395
396 def process_iop(directory, max_distance=400, filter_t=2):
397     """
398     Processes the entire directory using the file_to_data function.
399
400     Parameters
401     -----
402     directory (str): Path to the LiDAR text file.
403
404     Returns
405     -----
406     pandas.DataFrame: A DataFrame containing the processed LiDAR data with background
407     noise correction and
408     spatial coordinates (x, z, distance) for the entire directory.
409
410     Examples
411     -----
412
413     """
414     df = pd.DataFrame()
415     for file in os.listdir(directory):
416         if file.endswith('.txt'):
417             temp_df = file_to_data(os.path.join(directory, file), max_distance=
418             max_distance, filter_t=filter_t)
419             df = pd.concat((df, temp_df), ignore_index=True)
420
421     return df
422
423
424 def plot_contour_scan(scan_df, column="analog_rcs", title=None, x_limits=None, y_limits=
425     None, method='linear', surface=None, mark_max=False):
426     x = scan_df["x"]
427     z = scan_df["z"]
428     value = scan_df[column]
429
430     # Create a grid
431     xi = np.linspace(x.min(), x.max(), 100)
432     zi = np.linspace(z.min(), z.max(), 100)
433     Xi, Zi = np.meshgrid(xi, zi)
434
435     # Interpolate data
436     Ai = griddata((x, z), value, (Xi, Zi), method=method) # methods: linear, nearest,
437     cubic
438
439     # Plot contour map
440     plt.figure(figsize=(15, 6))
441     contour = plt.contourf(Xi, Zi, Ai, cmap="turbo", levels=30) # gist_ncar

```

```
430 plt.colorbar(label="Relative Backscatter")
431 plt.xlabel("Distance (m)")
432 plt.ylabel("Height (m)")
433 if title == None:
434     plt.title("Contour Map of Normalized Backscattter")
435 else:
436     plt.title(title)
437 if surface:
438     plt.plot(surface['x'], surface['z'])
439 if x_limits:
440     plt.xlim(x_limits[0], x_limits[1])
441 if y_limits:
442     plt.ylim(y_limits[0], y_limits[1])
443
444
445 if mark_max==True:
446     # find point of max
447     max_rcs_index = scan_df['analog_rcs'].idxmax()
448     max_backscatter_data = scan_df.loc[max_rcs_index]
449
450     max_x = max_backscatter_data['x']
451     max_z = max_backscatter_data['z']
452
453     plt.annotate(
454         'Maximum Backscatter',
455         xy=(max_x, max_z),
456         xytext=(100,25),
457         arrowprops=dict(
458             arrowstyle='->',
459             color='red',
460             lw=2
461         ),
462         fontsize=12,
463         color='green'
464     )
465
466 plt.show()
```