

▼ Initial Setup

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from datetime import datetime, timedelta
6 FILENAMES = ['participants', 'blood_pressure', 'heart_rate', 'hrv_measurements', 'surveys', 'scales_description']
7 URL = 'https://raw.githubusercontent.com/Welltory/hrv-covid19/master/data/'
8 EXTENSION = '.csv'
9 dfs = {}
10 for fn in FILENAMES:
11     dfs[fn] = pd.read_csv(URL + fn + EXTENSION)
```

▼ COVID Preprocessing and Feature Extraction

1. Maps specific survey responses to 0 or 1, 0 being low likelihood of COVID, 1 being likelihood of COVID
2. Envelopes responses to get an approximate start and stop time for symptoms.

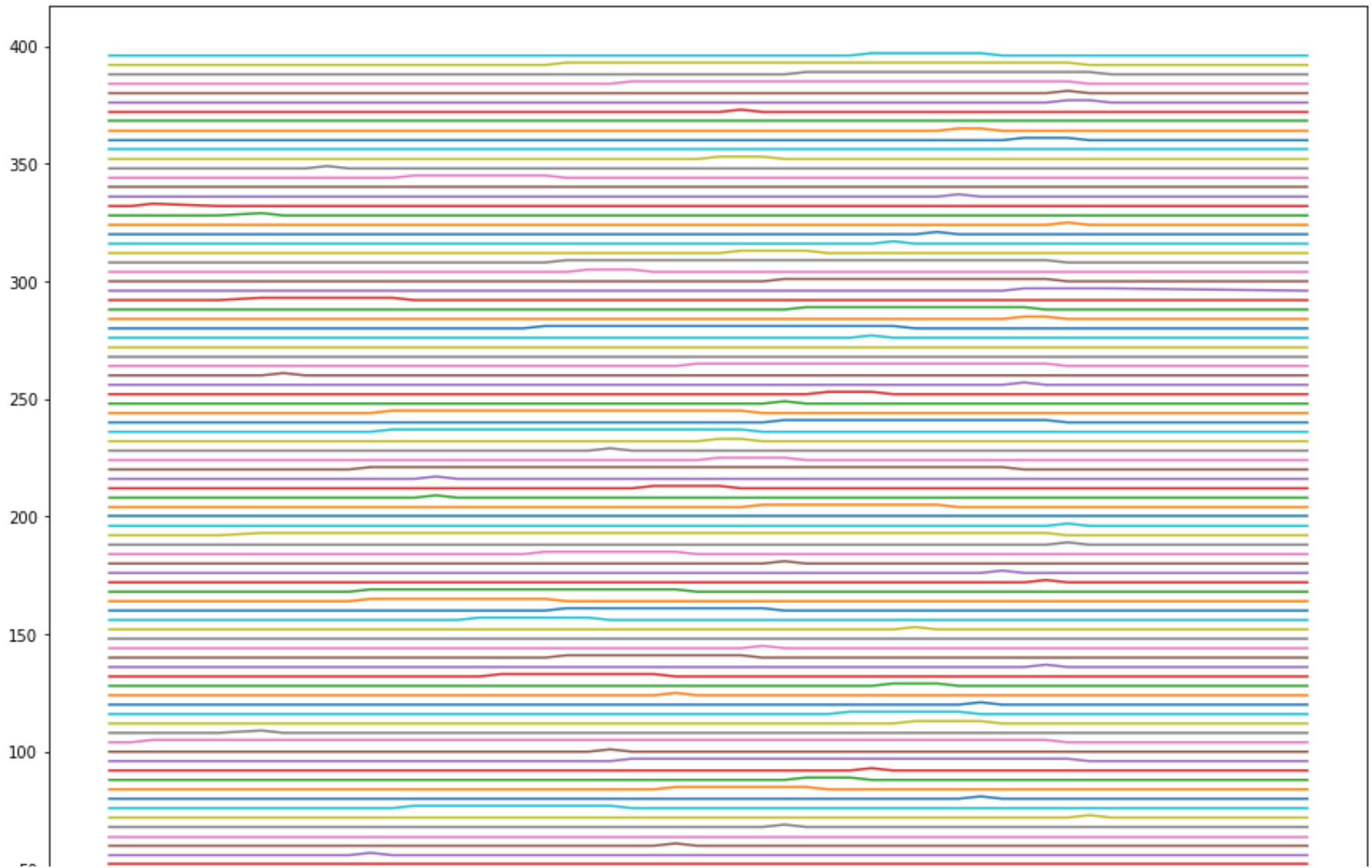
```
1 # Mapping for covid variables
2 # https://github.com/Welltory/hrv-covid19/blob/master/data/scales\_description.csv
3 vals = {'H': 1, 'M': 1, 'L': 0}
4 keys = ['S_CORONA', 'S_COVID_OVERALL']
5 maps = [{1: 'H', 2: 'H', 3: 'M', 4: 'L', 5: 'L'},
6         {1: 'L', 2: 'M', 3: 'M', 4: 'H', 5: 'H', 6: 'H'}]
7
8 # Initial processing
9 df = dfs['surveys'].copy()
10 df['created_at'] = pd.to_datetime(df['created_at'])
```

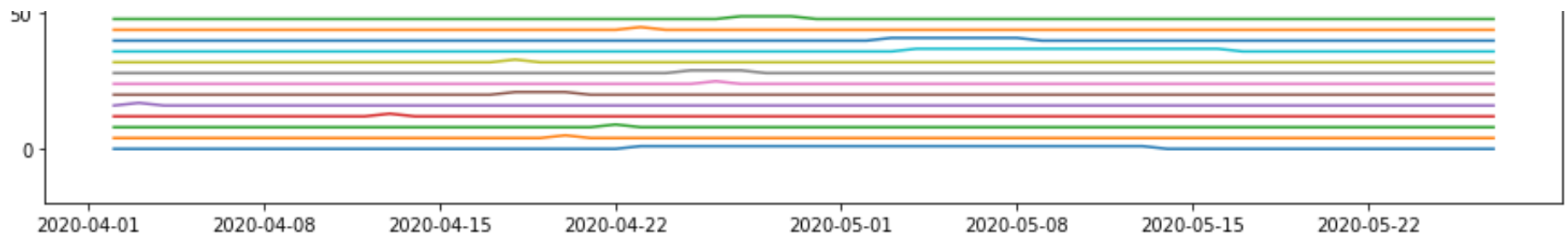
```

11 df = df.loc[df['scale'].isin(keys)]
12 for i, key in enumerate(keys):
13     df.loc[df['scale'] == key, 'value'] = df.loc[df['scale'] == key, 'value'].map(maps[i])
14 df = df.drop(columns=['text', 'scale'])
15 df = df.replace({'value': vals})
16
17 # Pivot to make index the userid and columns each day
18 df = df.copy().pivot_table(index='user_code', columns='created_at', values='value', aggfunc=np.max).fillna
19 df = df.reset_index()
20 df = df.set_index('user_code')
21 df.columns.name = None
22 df.index.name = None
23
24 # Pad an extra day for envelope
25 df.insert(0, min(df) - timedelta(1), 0.0)
26 df.insert(0, max(df) + timedelta(1), 0.0)
27
28 # Envelope covid data
29 def env(y):
30     y1 = y.replace(to_replace=0, method='ffill')
31     y2 = y.replace(to_replace=0, method='bfill')
32     y3 = y1 & y2
33     return y3
34 df = df.astype('int').apply(lambda y: env(y), axis=1)
35
36 # # Remove empty rows
37 # df = df[(df!=0).any(1)].copy()
38
39 # # Remove rows where covid symptoms are less than 1 week
40 # df = df[df.sum(axis=1) > 7].copy();
41
42 # Plot
43 # Each line is a survey response
44 plt.figure(figsize=(15,12))
45 for i in range(len(df)):
46     x = df.iloc[i,1:].index
47     y = df.iloc[i,1:].values + i*4

```

```
47 y = df.iloc[:,1].values + 104
48 plt.plot(x, y)
49
50 # Save for later
51 df_covid = df.copy()
52
53 #print(df_covid.head())
```





▼ HRV Preprocessing

Output: df_hrv_pp - a dictionary of {user_code : dataframe}

Each dataframe contains the HRV time series of that particular user

One random participant's HRV is plotted at the end

```

1 # HRV Data (Colman)
2
3 df_hrv = dfs['hrv_measurements']
4
5 # print the columns to see what data does it have
6 #print(df_hrv.columns)
7
8 # let's look at the datetime column
9 #print(df_hrv['measurement_datetime'].unique())
10
11 # We get a list of dates and use it as the x-axis of the time series (TODO: probably need a better way)
12 #dlist = pd.to_datetime(df_hrv['measurement_datetime']).dt.date
13 #dlist = pd.to_datetime(dlist).dt.normalize() # convert from object to datetime[ns]
14 #dlist = dlist.unique() # turns out only 172 days have measurements
15
16 # Note that the measurement is in date + time, but we are predicting the onset date only
17 # so we have to combine measurement of the same date into one
18

```

```

19 # get list of column names for new dataframes
20 #column_names = df_hrv.columns
21 #column_names = column_names.drop(labels=['user_code', 'measurement_datetime']) # remove unnecessary column
22
23 # prepare to construct a dataframe for each participant (TODO: is there better approach for a 3D dataframe)
24 df_hrv_pp = {}
25 # get list of participants
26 plist = df_hrv['user_code'].unique()
27
28 # start with an naive approach: we only get the first entry and drop the remaining ones
29 # (i.e., we take the first measurement of the day to represent the whole day)
30 # iterate participant
31 for p in plist:
32     # filter out measurements and drop NaN
33     tmp = df_hrv.loc[df_hrv['user_code'] == p].copy().dropna()
34
35     # TODO drop data if it does not fulfill requirement (e.g., less than 5 entries)
36     if tmp.shape[0] < 5:
37         continue
38
39     # convert datetime to date only
40     tmp['measurement_datetime'] = pd.to_datetime(tmp['measurement_datetime']).dt.date
41     tmp['measurement_datetime'] = pd.to_datetime(tmp['measurement_datetime']).dt.normalize() # convert from
42     #tmp['measurement_datetime'] = tmp['measurement_datetime'].dt.strftime('%Y-%m-%d')
43     # drop user_code column
44     tmp = tmp.drop(columns=['user_code', 'rr_code', 'time_of_day', 'how_feel', 'how_mood', 'how_sleep', 'rr_
45
46     # set date time as key
47     tmp.set_index('measurement_datetime', inplace=True)
48
49     # check if datetime is unique
50     if not tmp.index.is_unique:
51         # calculate mean for duplicates
52         tmp = tmp.reset_index().pivot_table(columns=["measurement_datetime"]).T
53
54     # drop duplicates

```

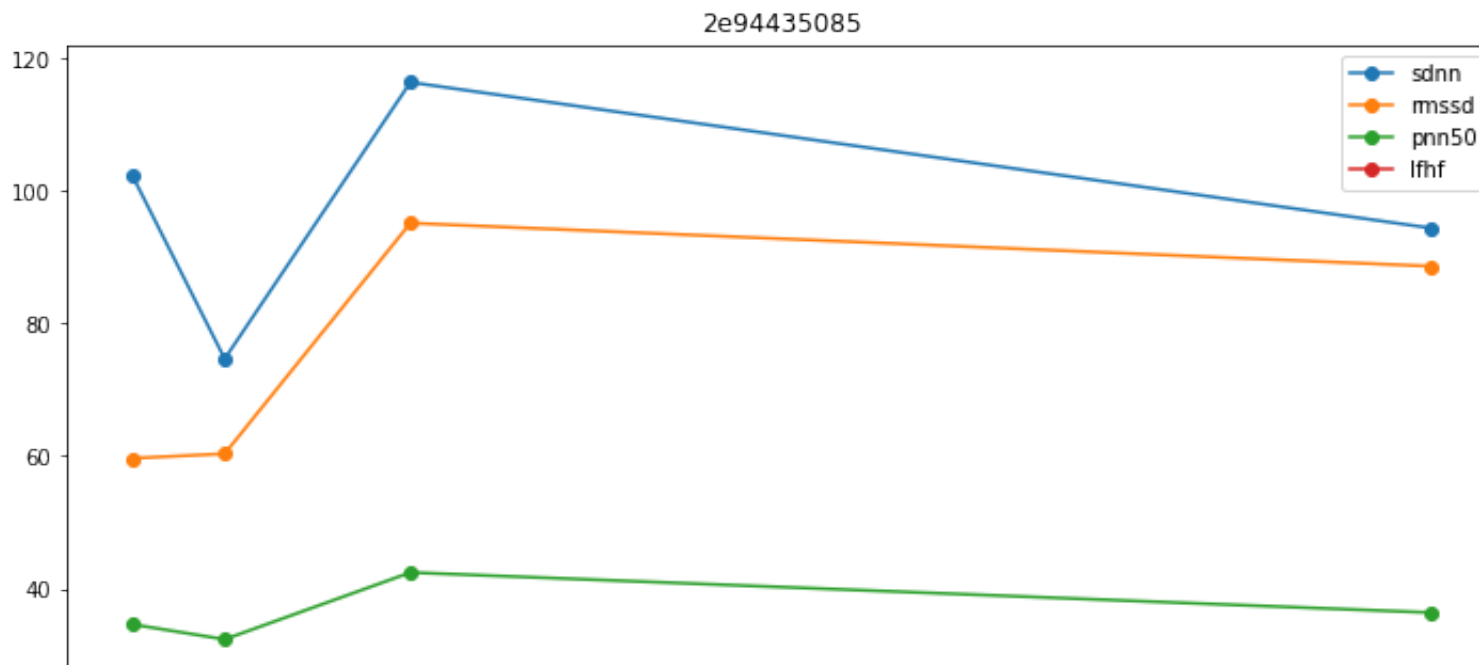
```

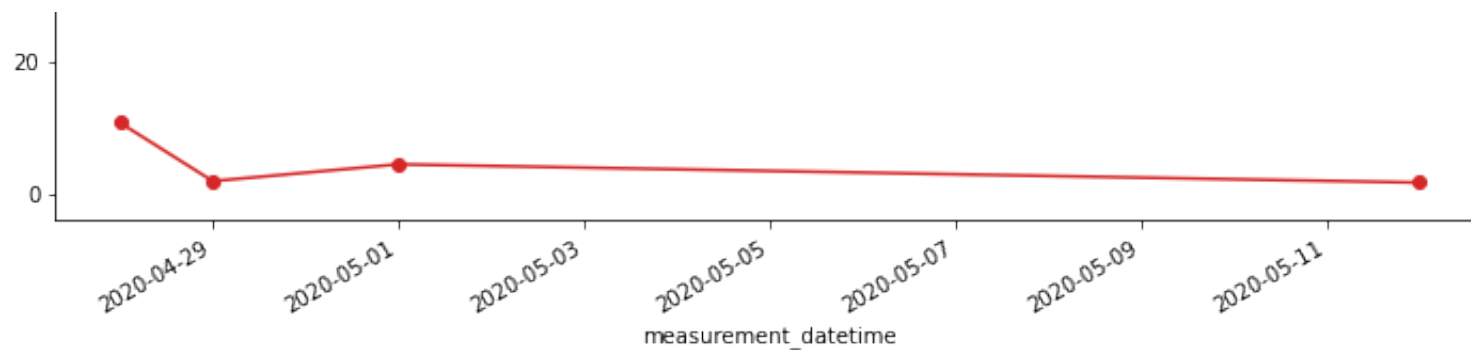
55 #tmp.drop_duplicates(keep='first')
56
57 # set dataframe
58 df_hrv_pp[p] = tmp
59
60 print("Extracted: " + str(len(df_hrv_pp)) + " participants")
61
62 # let's randomly plot one participant
63 import random
64 user_code, df_hrv_random = random.choice(list(df_hrv_pp.items()))
65
66 #print(df_hrv_random)
67 df_hrv_random.plot(y=['sdnn','rmssd','pnn50','lfhf'], kind='line', marker='o', figsize=(12, 8), title=user
68
69 # also, fcf3ea75b0 is interesting to look at
70 user_code = 'fcf3ea75b0'
71 df_hrv_pp[user_code].plot(y=['sdnn','rmssd','pnn50','lfhf'], kind='line', marker='o', figsize=(12, 8), tit

```

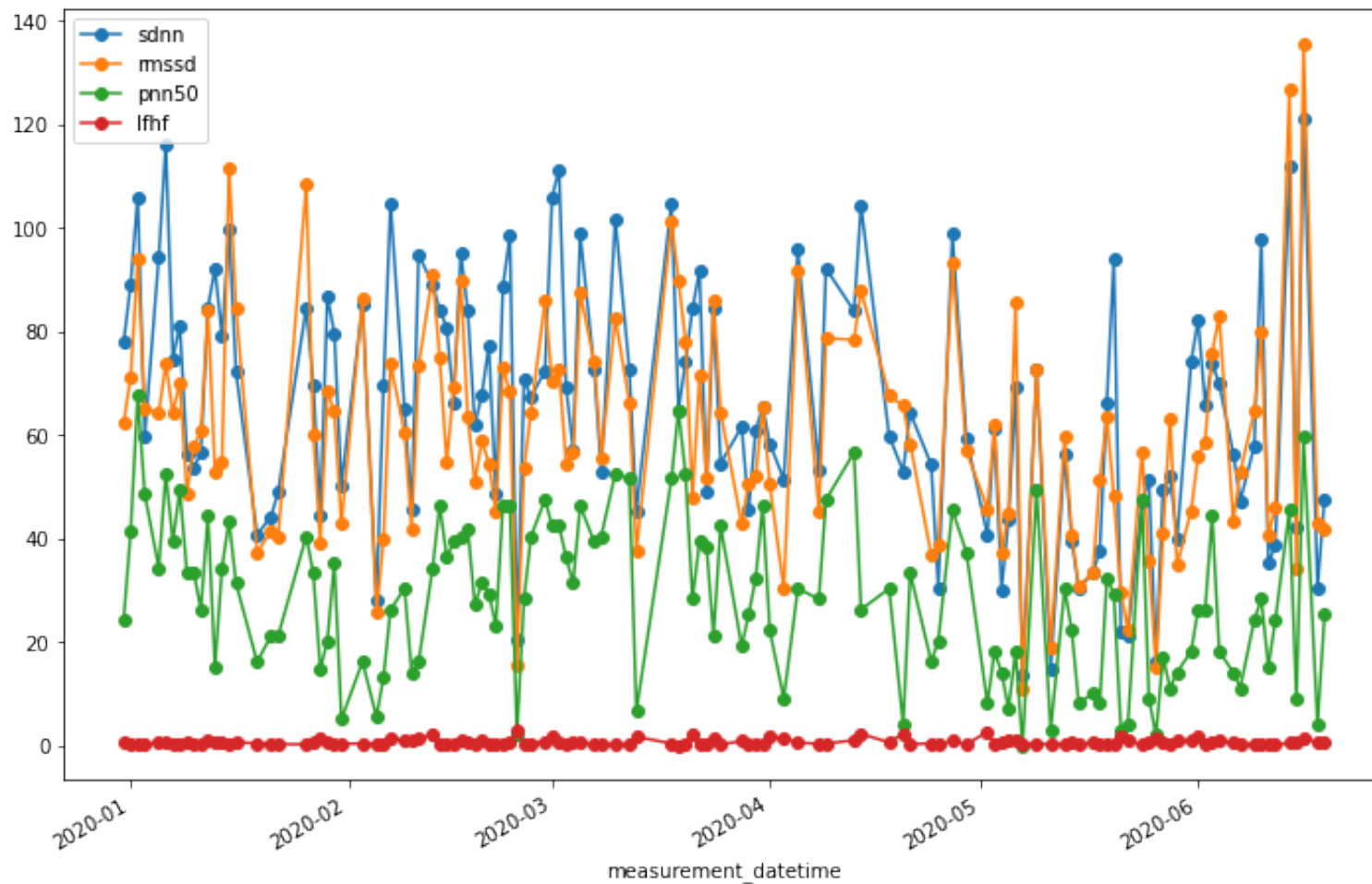
Extracted: 41 participants

<matplotlib.axes._subplots.AxesSubplot at 0x7f848f282350>





fcf3ea75b0



▼ Sleep Preprocessing

```
1 # Sleep data (Hunter)
2 # 1. Bring in sleep data
3 # 2. Filter by patient ID and dates
4 # 3. Prepare dataframe for LR or other ML model
5
6 #print(df_covid)
7
8 old_df_sleep = dfs['sleep']
9 # drop columns that aren't filled for all users
10 old_df_sleep = old_df_sleep.drop(columns=['sleep_awake_duration', 'sleep_rem_duration', 'sleep_light_durat
11 old_df_sleep = old_df_sleep.drop(columns=['pulse_min', 'pulse_max', 'pulse_average'])
12 old_df_sleep = old_df_sleep.drop(columns=['sleep_begin', 'sleep_end'])
13 # assign 1 if user_code is in df_p, 0 if not
14 old_df_sleep = old_df_sleep.assign(InDFP = old_df_sleep.user_code.isin(df_covid.index).astype(int))
15 # only keep user_codes that were in df_p, drop extra column
16 old_df_sleep = old_df_sleep[old_df_sleep.InDFP != 0].drop(columns=['InDFP'])
17 # print(old_df_sleep)
18
19 sleep_columns = list(old_df_sleep.day.unique())
20 sleep_index = list(old_df_sleep.user_code.unique())
21 user_dict = {}
22 for code in sleep_index:
23     user_dict[code] = {}
24
25 sleep_user_code = old_df_sleep.user_code.tolist()
26 sleep_day = old_df_sleep.day.tolist()
27 sleep_duration = old_df_sleep.sleep_duration.tolist()
28
29 for i in range(0, len(sleep_user_code)):
30     user_dict[sleep_user_code[i]][sleep_day[i]] = sleep_duration[i]
31
```



```
32 #print(user_dict)
33
34 sleep_dict = {}
35 for date in sleep_columns:
36     sleep_dict[date] = []
37     for code in sleep_index:
38         if date not in user_dict[code].keys():
39             sleep_dict[date].append(0)
40         else:
41             sleep_dict[date].append(user_dict[code][date])
42
43 df_sleep = pd.DataFrame(sleep_dict, index = sleep_index)
44
45 print(df_sleep)
```

	2020-01-21	2020-01-30	2020-01-31	2020-02-02	2020-02-03	\
276ab22485	9543.0	0.0	0.0	0.0	0.0	
4985083f4d	0.0	29265.0	24771.0	11410.0	31705.0	
6be5033971	29400.0	29700.0	35100.0	30300.0	35700.0	
9871ee5e7b	0.0	0.0	0.0	0.0	0.0	
a1c2e6b2eb	0.0	0.0	0.0	0.0	0.0	
c174f32d88	0.0	0.0	0.0	0.0	0.0	
fcf3ea75b0	0.0	0.0	0.0	0.0	0.0	
	2020-02-05	2020-02-09	2020-02-12	2020-02-13	2020-02-14	... \
276ab22485	0.0	0.0	0.0	0.0	0.0	...
4985083f4d	24492.0	14522.0	5130.0	5235.0	29883.0	...
6be5033971	35700.0	33300.0	36600.0	36300.0	33300.0	...
9871ee5e7b	0.0	0.0	0.0	0.0	0.0	...
a1c2e6b2eb	0.0	0.0	0.0	0.0	0.0	...
c174f32d88	0.0	0.0	0.0	0.0	0.0	...
fcf3ea75b0	0.0	0.0	0.0	0.0	0.0	...
	2020-04-13	2020-04-14	2020-04-15	2020-04-16	2020-04-17	\
276ab22485	0.0	0.0	0.0	0.0	0.0	
4985083f4d	0.0	0.0	0.0	0.0	0.0	
6be5033971	0.0	0.0	0.0	0.0	0.0	
9871ee5e7b	0.0	0.0	0.0	0.0	0.0	
a1c2e6b2eb	0.0	0.0	0.0	0.0	0.0	
c174f32d88	35940.0	40050.0	35580.0	42870.0	34980.0	
fcf3ea75b0	0.0	0.0	0.0	0.0	0.0	
	2020-04-19	2020-04-20	2020-04-21	2020-04-22	2020-04-23	
276ab22485	0.0	0.0	0.0	0.0	0.0	
4985083f4d	0.0	0.0	0.0	0.0	0.0	
6be5033971	0.0	0.0	0.0	0.0	0.0	
9871ee5e7b	0.0	0.0	0.0	0.0	0.0	
a1c2e6b2eb	0.0	0.0	0.0	0.0	0.0	
c174f32d88	47460.0	41610.0	37170.0	39720.0	37710.0	
fcf3ea75b0	0.0	0.0	0.0	24841.0	30224.0	

[7 rows x 135 columns]

▼ Wearable Preprocessing

```
1 #Wearable Data
2 #TODO:
3 # 1 : Load in wearable data
4 # 2 : Filter by patient ID and dates
5 # 3 : Prepare dataframe for LR or other ML model
6
7 df_wearable = dfs['wearables'].copy()
8
9 #drop columns that aren't filled for most users
10 #adjust if necessary later
11 df_wearable = df_wearable.drop(columns=['resting_pulse', 'average_spo2_value', 'body_temperature_avg', 'st
12 # assign 1 if user_code is in df_p, 0 if not
13 df_wearable = df_wearable.assign(InDFP = old_df_sleep.user_code.isin(df_covid.index).astype(int))
14 df_wearable = df_wearable[df_wearable['user_code'].notna()]
15
16 # only keep user_codes that were in df_p, drop extra column
17 df_wearable = df_wearable[df_wearable.InDFP != 0].drop(columns=['InDFP'])
18
19 #remove user codes which are NaN
20 #Note: Check why this happened
21 df_wearable = df_wearable[df_wearable['user_code'].notna()]
22
23 df_wearable = df_wearable[df_wearable['user_code'].notna()]
24
25 df_wearable['date'] = df_wearable.apply(lambda _: '', axis=1)
26
27 df_wearable_cp= df_wearable.copy()
28 df_wearable_cp = df_wearable_cp.dropna()
29 # df_wearable_cp2= df_wearable.copy()
30
31 # wearable_columns = list(df_wearable.day.unique())
```

```

32 # wearable_index = list(df_wearable.user_code.unique())
33 # user_dict = {}
34 # for code in wearable_index:
35 #     user_dict[code] = {}
36
37 # wearable_user_code = df_wearable.user_code.tolist()
38 # wearable_day = df_wearable.day.tolist()
39 # wearable_duration = df_wearable.sleep_duration.tolist()
40 # print(df_wearable_cp.columns)
41 # print(df_wearable.iloc[0][1])
42
43 #Filter Dates
44 # for i, r in enumerate(df_wearable['day']):
45 #     str_list = r.split("-")
46 #     for j in range(0, len(str_list)):
47 #         if str_list[j][0] == '0':
48 #             str_list[j] = str_list[j][1:]
49 #     new_str_list = [str_list[1], str_list[2], str_list[0]]
50 #     df_wearable.loc[i, 'day'] = '/'.join(new_str_list)
51
52
53
54
55
56 # df_wearable_mapping = df_wearable.assign(StartDate = df_wearable.day.isin(df_p.symptoms_onset).astype(int))
57 # df_wearable_mapping = df_wearable_mapping.dropna()
58
59 # print(df_wearable_mapping)
60 # print(len(list(df_wearable_mapping.user_code.unique()))))
61
62 df_wearable_cp['combo_index'] = df_wearable_cp['user_code'].str.cat(df_wearable_cp['day'], sep = "_")
63 # df_wearable_cp = df_wearable_cp.drop(columns=['user_code', 'day'])
64 # df_wearable_cp = df_wearable_cp.set_index('combo_index')
65 print(df_wearable_cp.shape)

```

(910, 11)

▼ Classification

▼ HRV Classification

```
1 # # assign 1 if user_code is in df_p, 0 if not
2 # old_df_sleep = old_df_sleep.assign(InDFP = old_df_sleep.user_code.isin(df_covid.index).astype(int))
3 # # only keep user_codes that were in df_p, drop extra column
4 # old_df_sleep = old_df_sleep[old_df_sleep.InDFP != 0].drop(columns=['InDFP'])
5
6 #temp_df_hrv = pd.DataFrame(df_hrv_pp)
7 sklearn_covid = df_covid.assign(InDF = df_covid.index.isin(df_hrv_pp.keys()).astype(int))
8 sklearn_covid = sklearn_covid[sklearn_covid.InDF != 0].drop(columns=['InDF'])
9 #print(sklearn_covid.shape)
10 covid_codes = sklearn_covid.index.tolist()
11 covid_dates = sklearn_covid.columns.tolist()
12
13 temp_df_hrv = df_hrv_pp.copy()
14 for key in df_hrv_pp.keys():
15     if key not in sklearn_covid.index:
16         del temp_df_hrv[key]
17
18 col_date = 'measurement_datetime'
19 dates = []
20 all_dates = []
21 for key in temp_df_hrv.keys():
22     these_dates = temp_df_hrv[key].index.tolist()
23     for date in these_dates:
24         if date not in dates and date in covid_dates:
25             dates.append(date)
26 # print(dates)
27 for i in range(len(dates)):
```

```

28     dates[i] = dates[i].to_pydatetime()
29     dates[i] = dates[i].strftime("%Y-%m-%d")
30 # print(dates)
31 # print(len(temp_df_hrv.keys()), len(dates))
32 # print(temp_df_hrv[covid_codes[0]])
33 # print(temp_df_hrv[covid_codes[0]].loc[dates[0]]['amo'])
34
35 combo_index = []
36 for code in covid_codes:
37     for date in dates:
38         combo_index.append(code + "_" + str(date))
39
40 full_hrv_dict = {}
41 full_covid_dict = {}
42 hrv_columns = temp_df_hrv[covid_codes[0]].columns.tolist()
43 for col in hrv_columns:
44     full_hrv_dict[col] = []
45 full_hrv_dict['covid'] = []
46
47 errors = 0
48 for code in covid_codes:
49     for date in dates:
50         for col in hrv_columns:
51             try:
52                 full_hrv_dict[col].append(temp_df_hrv[code].loc[date][col])
53             except:
54                 full_hrv_dict[col].append(0.0)
55         full_hrv_dict['covid'].append(int(sklearn_covid.loc[code][date]))
56
57 sklearn_hrv = pd.DataFrame(full_hrv_dict, index=combo_index)
58 # test_clf = LogisticRegression(random_state=0).fit(df_covid, df_hrv_pp)

```

```
1 print(sklearn_hrv.shape)
2 # print(sklearn_hrv.head)
3 sklearn_hrv = sklearn_hrv[sklearn_hrv.amo != 0.0]
4 print(sklearn_hrv.shape)
```

```
(1472, 14)
```

```
(328, 14)
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.svm import SVC, LinearSVC, NuSVC
4 from sklearn.neural_network import MLPClassifier
```

```

1 X_train, X_test, y_train, y_test = train_test_split(sklearn_hrv.drop(["covid"], axis=1), sklearn_hrv["covi
2
3 lr = LogisticRegression(random_state=42, max_iter=1000).fit(X_train, y_train)
4 lr_score = lr.score(X_test, y_test)
5 print("Logistic Regression: " + str(lr_score))
6
7 mlpc = MLPClassifier(random_state=42, max_iter=1000).fit(X_train, y_train)
8 mlpc_score = mlpc.score(X_test, y_test)
9 print("MultiLayer Perceptron: " + str(mlpc_score))
10
11 # SVMs are SLOW but they eventually complete lol
12 linsvc = LinearSVC(random_state=42).fit(X_train, y_train)
13 linsvc_score = linsvc.score(X_test, y_test)
14 print("SVM-LinearSVC: " + str(linsvc_score))
15
16 nusvc = NuSVC(random_state=42).fit(X_train, y_train)
17 nusvc_score = nusvc.score(X_test, y_test)
18 print("SVM-NuSVC: " + str(nusvc_score))
19
20 svclin = SVC(kernel="linear", random_state=42).fit(X_train, y_train)
21 svclin_score = svclin.score(X_test, y_test)
22 print("SVM-SVC (linear kernel): " + str(svclin_score))
23
24 svcrbf = SVC(kernel="rbf", random_state=42).fit(X_train, y_train)
25 svcrbf_score = svcrbf.score(X_test, y_test)
26 print("SVM-SVC (rbf kernel): " + str(svcrbf_score))

```

Logistic Regression: 0.6585365853658537

MultiLayer Perceptron: 0.573170731707317

/usr/local/lib/python3.7/dist-packages/sklearn/svm/_base.py:1208: ConvergenceWarning: Liblinear failed to converge, raise ConvergenceWarning,

SVM-LinearSVC: 0.5609756097560976

SVM-NuSVC: 0.5853658536585366

SVM-SVC (linear kernel): 0.7073170731707317

SVM-SVC (rbf kernel): 0.5487804878048781

▼ Wearable Classification

```
1 # ['user_code', 'day', 'pulse_average', 'pulse_min', 'pulse_max',
2 #     'steps_count', 'steps_speed', 'basal_calories_burned',
3 #     'total_calories_burned', 'date']
4 print(df_wearable_cp.shape)
5 wearable_index = df_wearable_cp.index.tolist()
6 bad_wearable_index_i = []
7 for i in wearable_index:
8     if df_wearable_cp['combo_index'][i] not in combo_index:
9         bad_wearable_index_i.append(i)
10
11 sklearn_wearable = df_wearable_cp.drop(labels=bad_wearable_index_i, axis=0)
12 sklearn_wearable = sklearn_wearable.set_index('combo_index')
13 sklearn_wearable['covid'] = ""
14
15 for index in sklearn_wearable.index:
16     if sklearn_wearable.loc[index][0] in sklearn_covid.index and sklearn_wearable.loc[index][1] in dates:
17         sklearn_wearable.at[index, 'covid'] = int(sklearn_covid.loc[sklearn_wearable.loc[index][0]][sklearn_covid.columns[1]])
18
19 sklearn_wearable = sklearn_wearable.drop(columns=['user_code', 'day'])
20 sklearn_wearable = sklearn_wearable.apply(pd.to_numeric)
21 sklearn_wearable = sklearn_wearable.drop(columns=['date'])
22 print(sklearn_wearable.shape)
23 # print(sklearn_wearable.head())
```

(910, 11)
(204, 8)

```

1 X_train, X_test, y_train, y_test = train_test_split(sklearn_wearable.drop(["covid"], axis=1), sklearn_wear
2
3 lr = LogisticRegression(random_state=42, max_iter=1000).fit(X_train, y_train)
4 lr_score = lr.score(X_test, y_test)
5 print("Logistic Regression: " + str(lr_score))
6
7 mlpc = MLPClassifier(random_state=42, max_iter=1000).fit(X_train, y_train)
8 mlpc_score = mlpc.score(X_test, y_test)
9 print("MultiLayer Perceptron: " + str(mlpc_score))
10
11 # SVMs are SLOW but they eventually complete lol
12 linsvc = LinearSVC(random_state=42).fit(X_train, y_train)
13 linsvc_score = linsvc.score(X_test, y_test)
14 print("SVM-LinearSVC: " + str(linsvc_score))
15
16 nusvc = NuSVC(random_state=42).fit(X_train, y_train)
17 nusvc_score = nusvc.score(X_test, y_test)
18 print("SVM-NuSVC: " + str(nusvc_score))
19
20 svclin = SVC(kernel="linear", random_state=42).fit(X_train, y_train)
21 svclin_score = svclin.score(X_test, y_test)
22 print("SVM-SVC (linear kernel): " + str(svclin_score))
23
24 svcrbf = SVC(kernel="rbf", random_state=42).fit(X_train, y_train)
25 svcrbf_score = svcrbf.score(X_test, y_test)
26 print("SVM-SVC (rbf kernel): " + str(svcrbf_score))

```

Logistic Regression: 0.6274509803921569

MultiLayer Perceptron: 0.45098039215686275

SVM-LinearSVC: 0.49019607843137253

SVM-NuSVC: 0.6862745098039216

/usr/local/lib/python3.7/dist-packages/sklearn/svm/_base.py:1208: ConvergenceWarning: Liblinear failed to converge, increasing max_iter may resolve this.

SVM-SVC (linear kernel): 0.5098039215686274

SVM-SVC (rbf kernel): 0.5098039215686274

▼ Other

```
1 # boundary check (make sure the input is of length 10 and within the acceptable range)
2 # note that we assumed all inputs are integers
3 def checkBoundary(scores):
4     # check length
5     if len(scores) != 10:
6         return False
7     # check boundary
8     #return (scores[0] < 1 or scores[1] < 1 or scores[2] < 1 or scores[3] < 1 or scores[4] < 1 or scores[5]
9     #scores[0] > 6 or scores[1] > 5 or scores[2] > 5 or scores[3] > 5 or scores[4] > 5 or scores[5] > 5 or sco
10    return scores[0] < 1 or scores[0] > 6
11
12 # return true if user is negative (no symptoms)
13 def checkNegative(scores):
14     #return scores[0] == 1 and scores[1] >= 3 and scores[1] <= 5 and scores[2] == 1 and scores[3] == 1 and s
15     return scores[0] == 1
16
17 # input first time report of symptoms + symptom score [overall, corona, cough, fever, breath, fatigue, pai
18 # output NOGCV labeled as 0 = 4 (negative, onset, general, critical, ventilator)
19 # output None if data is invalid
20 # right now this is simplified to only consider S_COVID_OVERALL
21 def covidSymptomScore(reportedFirstTime, scores):
22     # boundary check
23     if checkBoundary(scores):
24         print("Invalid scores (out of bound)")
25         return None
26
27     # set score for each category
28     result = [0] * 10
29
30     # if already has covid (reported already)
31     # report g, c, or v depending on severeness of symptoms
```

```
32  if reportedFirstTime:
33      if checkNegative():
34          return 0
35      if scores[0] <= 4:
36          return 2
37      elif scores[0] == 5:
38          return 3
39      elif scores[0] == 6:
40          return 4
41      else:
42          print("Invalid scores (unknown)")
43          return None
44
45  # doesnt get covid yet
46  else:
47      if checkNegative():
48          return 0 # still negative
49      else:
50          return 1 # tested positive, marked as onset (1)
```

Preprocessing data using the flow mentioned in "Pre-Emption of Affliction Severity Using HRV Measurements from a Smart Wearable Case-Study on SARS-Cov-2 Symptoms"

<https://www.mdpi.com/1424-8220/20/24/7068>

```

1 # get participants data
2 df_p = dfs['participants'].dropna()
3 df_p = df_p.drop(columns=['city', 'country']) # drop irrelevant columns
4
5 # map gender and age group to a number (note that if we want to use it as a feature, it will be better to
6 dict_gender = {'m':0, 'f':1}
7 dict_age = {'18-24':0, '25-34':1, '35-44':2, '45-54':3, '55-64':4, '65-74':5}
8 df_p['gender'] = df_p['gender'].map(dict_gender)
9 df_p['age_range'] = df_p['age_range'].map(dict_age)
10
11
12
13 print(df_p)

```

	user_code	gender	age_range	height	weight	symptoms_onset
1	013f6d3e5b	1	0	174.00	77.300	5/15/2020
2	01bad5a519	0	3	178.00	92.000	4/5/2020
3	0210b20eea	1	1	169.00	60.000	5/6/2020
4	024719e7da	1	3	158.00	68.500	5/27/2020
9	0bdfbddb2b	1	0	159.00	73.500	4/1/2020
..
178	f9edcb7056	1	5	154.94	130.300	3/16/2020
179	fcf3ea75b0	1	3	168.00	92.644	5/1/2020
180	fd387f6269	1	2	165.00	115.439	5/1/2020
182	fde84801d8	1	3	168.00	79.500	4/16/2020
184	fe6c1b1349	1	1	173.00	53.000	5/3/2020

[136 rows x 6 columns]

