

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Mon Dec 23 16:32:22 2019
4
5  @author: Alan.Toppen
6  """
7
8  import pandas as pd
9  import numpy as np
10 import sqlalchemy as sq
11 import io
12 import boto3
13 from multiprocessing import get_context, Pool
14 import itertools
15 import re
16 import os
17 import sys
18 from datetime import datetime, timedelta
19 import time
20 import random
21
22 s3 = boto3.client('s3')
23 events_bucket = 'gdot-spm'
24 config_bucket = events_bucket
25
26
27 def read_atspm_query(query):
28     engine = sq.create_engine('mssql+pyodbc://atspm',
29                               pool_size=20)
30
31     with engine.connect() as con:
32         df = pd.read_sql_query(query, con=con)
33     return df
34
35
36 def get_eventlog_data_db(signalid, date_str):
37
38     start_date = date_str # date_.strftime('%Y-%m-%d %H:%M:%S.%f')[:-5]
39     end_date = (pd.Timestamp(date_str) + pd.DateOffset(days=1) - pd.DateOffset(seconds=0.1))\
40               .strftime('%Y-%m-%d %H:%M:%S.%f')[:-5]
41
42     df = read_atspm_query("""
43         SELECT * FROM Controller_Event_Log
44         WHERE SignalID = '{}'
45         AND Timestamp BETWEEN '{}' AND '{}'
46         ORDER BY SignalID, Timestamp, EventCode, EventParam
47         """.format(signalid.zfill(5),
48                   start_date,
49                   end_date))
50     return df
51
52
53 def get_eventlog_data(bucket, signalid, dates):
54     for date_ in dates:
55         date_str = date_.strftime('%F')
56         df = pd.read_parquet('s3://{b}/atspm/date={d}/atspm_{s}_{d}.parquet'.format(
57             b=bucket, d=date_str, s=signalid))
58         df.Timestamp = df.Timestamp.dt.tz_localize(None)
59         df.SignalID = df.SignalID.astype('str')
60         df['date'] = date_
61         yield df
62
63
64
65

```

```

66 def get_det_config(bucket, date_, leading_zeros=False):
67     date_str = date_.strftime('%F')
68     objs = s3.list_objects(Bucket=bucket, Prefix=f'config/atspm_det_config_good/date={date_str}/')
69     keys = [obj['Key'] for obj in objs['Contents']]
70
71     def f(bucket, key):
72         with io.BytesIO() as data:
73             s3.download_fileobj(
74                 Bucket=bucket,
75                 Key=key,
76                 Fileobj=data)
77
78             dc = pd.read_feather(data)\
79                 .assign(SignalID = lambda x: x.SignalID.astype('str'))\
80                 .assign(Detector = lambda x: x.Detector.astype('int64'))\
81                 .reset_index(drop=True)
82
83             dc['date'] = date_
84
85             if leading_zeros:
86                 dc['SignalID'] = dc['SignalID'].str.zfill(5)
87
88             return dc
89
90     return pd.concat(map(lambda k: f(bucket, k), keys))
91     #.rename(columns={'CallPhase': 'Call Phase'})\
92
93
94 def get_det_configs(bucket, dates, leading_zeros=False):
95     return pd.concat([get_det_config(bucket, date_, leading_zeros) for date_ in dates])
96
97
98 def get_det_config_local(filename):
99
100     dc = pd.read_feather(filename)\
101         .assign(SignalID = lambda x: x.SignalID.astype('str'))\
102         .assign(Detector = lambda x: x.Detector.astype('int64'))\
103         .reset_index(drop=True)
104     #.rename(columns={'CallPhase': 'Call Phase'})\
105     return dc
106
107
108 def get_det_config_future(bucket, date_str):
109     key = 's3://{b}/atspm_det_config_good/date={d}/ATSPM_Det_Config_Good_Ozark.parquet'.format(
110         b=bucket, d=date_str)
111     print(key)
112     dc = pd.read_parquet(key).reset_index(drop=True)
113     return dc
114
115
116 # Works. Doesn't copy down. Use this for new grouping variable and copy_down to apply value across and down
117 def create_new_grouping_field(df, eventcodes, grouping_field, new_grouping_field, transform_func = lambda x: x):
118
119     if type(eventcodes) is list:
120         mask = df.EventCode.isin(eventcodes)
121     else:
122         eventcode = eventcodes
123         mask = df.EventCode==eventcode
124
125     df.loc[mask, new_grouping_field] = df.loc[mask, grouping_field].apply(transform_func)
126
127     return df
128
129
130

```

```

131 # Works. Two-step create new field and copy down. May need just a copy down.
132 def copy_updown(
133     df, eventcodes, new_field_name, group_fields, copy_field,
134     off_eventcode=None, direction='down', apply_to_timestamp='all'):
135     """
136     df - eventlog dataframe
137     eventcodes - EventCode(s) signifying event(s) to carry forward to subsequent events, e.g., 0 for PhaseStart
138     new_field_name - name of Event corresponding to EventCode, e.g., PhaseStart
139     group_fields - grouping(s) to which eventcode applies, e.g., [SignalID, EventParam] (Phase) for PhaseStart
140     copy_field - field identifying eventcode, e.g., Timestamp for PhaseStart
141     off_eventcode - optional value for where to stop copying up or down, otherwise goes to next value in
142     eventcodes
143     direction - 'up' for copy up, 'down' for copy down new_field_name
144     apply_to_timestamp - 'all' to fill all rows with timestamps of the eventcodes before copying up or down,
145         Example would be 31-Barrier which should renew with all events at that same
146         timestamp,
147         of which there are many starts and ends to phase intervals
148         'group' to fill all rows at the timestamp in the group
149         Example would be Recorded Split
150         None to not fill all rows at the timestamp.
151         Example detector off (81) or call off (44) events
152     """
153     if type(eventcodes) is list:
154         if sum(df.EventCode.isin(eventcodes)) == 0:
155             #print('Event Codes {} not in data frame'.format(','.join(map(str, eventcodes))))
156             return df
157         else:
158             df.loc[df.EventCode.isin(eventcodes), new_field_name] = df.loc[df.EventCode.isin(eventcodes),
159                 copy_field]
160     else:
161         eventcode = eventcodes
162         if sum(df.EventCode==eventcode) == 0:
163             #print('Event Code {} not in data frame'.format(eventcode))
164             return df
165         else:
166             df.loc[df.EventCode==eventcode, new_field_name] = df.loc[df.EventCode==eventcode, copy_field]
167
168     if apply_to_timestamp=='all':
169         df[new_field_name] = df.groupby(['SignalID', 'Timestamp'], group_keys=False)[new_field_name].transform(
170             'max') ## This seems to work
171     elif apply_to_timestamp=='group':
172         group_vars = list(set(['SignalID', 'Timestamp'] + group_fields))
173         df[new_field_name] = df.groupby(group_vars, group_keys=False)[new_field_name].transform('max') ## This
174         seems to work
175
176     if off_eventcode is not None:
177         df.loc[df.EventCode==off_eventcode, new_field_name] = -1
178
179     if direction == 'down':
180         df[new_field_name] = df.groupby(group_fields)[new_field_name].ffill()
181     elif direction == 'up':
182         df[new_field_name] = df.groupby(group_fields)[new_field_name].bfill()
183
184     if off_eventcode is not None:
185         df.loc[df[new_field_name]==-1, new_field_name] = None
186
187     return df
188
189
190

```

```

191 def copy_down(
192     df, eventcodes, new_field_name, group_fields, copy_field,
193     off_eventcode=None,
194     apply_to_timestamp='all'):
195     return copy_updown(
196         df, eventcodes, new_field_name, group_fields, copy_field,
197         off_eventcode=off_eventcode,
198         apply_to_timestamp=apply_to_timestamp,
199         direction='down')
200
201
202 def copy_up(
203     df, eventcodes, new_field_name, group_fields, copy_field,
204     off_eventcode=None,
205     apply_to_timestamp='all'):
206     return copy_updown(
207         df, eventcodes, new_field_name, group_fields, copy_field,
208         off_eventcode=off_eventcode,
209         apply_to_timestamp=apply_to_timestamp,
210         direction='up')
211
212
213 def widen(s, date_, det_config=None, source='s3'): # or source = 'db'
214
215     signalid = s
216     date0_ = date_ - pd.Timedelta(1, unit='D')
217     # date0_str = date0_.strftime('%F')
218     date_str = date_.strftime('%F')
219
220     print('{} | {} started.'.format(date_str, s))
221
222     if det_config is None:
223         det_config = get_det_configs(config_bucket, [date0_, date_])
224
225     dc = det_config[['SignalID', 'Detector', 'CallPhase', 'TimeFromStopBar', 'date']]
226     dc = dc[dc.SignalID==s]
227
228     if source=='s3':
229         df = pd.concat(get_eventlog_data(events_bucket, signalid, [date0_, date_]))
230     elif source=='db':
231         df = get_eventlog_data_db(signalid, date_str)
232
233     df = df[~df.EventCode.isin([43, 44])]
234
235     df = df.rename(columns={'TimeStamp': 'Timestamp'})
236     df = df.sort_values(['SignalID', 'Timestamp', 'EventCode', 'EventParam']).reset_index(drop=True)
237
238     print('{} | {} data queried from database.'.format(date_str, s))
239
240     # Map Detectors to Phases.
241     # Replace EventParam with Phase to align with other event types
242     # Add new field for Detector from original EventParam field
243     detector_codes = list(range(81,89))
244     dc2 = pd.concat([dc.assign(EventCode=d).rename(columns={'Detector': 'EventParam'}) for d in detector_codes])
245
246     df = pd.merge(
247         left=df,
248         right=dc2,
249         on=['SignalID', 'EventCode', 'EventParam', 'date'],
250         how='left')\
251         .reset_index(drop=True)
252
253     # Adjust Timestamp for detectors by adding TimeFromStopBar
254     df.Timestamp = df.Timestamp + pd.to_timedelta(df.TimeFromStopBar.fillna(0), 's')
255     df = df.sort_values(['SignalID', 'Timestamp', 'EventCode', 'EventParam'])

```

```

256
257 codes = df.EventCode.drop_duplicates().values.tolist()
258
259 # Rename Detector, Phase columns
260 df.loc[df.EventCode.isin(detector_codes), 'Detector'] = df.loc[df.EventCode.isin(detector_codes),
    'EventParam']
261 df.loc[df.EventCode.isin(detector_codes), 'Phase'] = df.loc[df.EventCode.isin(detector_codes), 'CallPhase']
262 df = df.drop(columns=['CallPhase', 'TimeFromStopBar'])
263
264 df = create_new_grouping_field(df, list(range(83,89)), ['SignalID', 'Detector'], 'DetectorFault')
265 df = copy_down(df, list(range(84,89)), 'DetectorFault', ['SignalID', 'Detector'], 'EventParam', off_eventcode=
    83)
266
267 ped_input_codes = [89, 90]
268 df.loc[df.EventCode.isin(ped_input_codes), 'PedInput'] = df.loc[df.EventCode.isin(ped_input_codes),
    'EventParam']
269
270 # Global (Signal-wide) copy-downs. Uses two-step function to create new field and copy down
271 if 31 in codes:
272     df = copy_down(df, 31, 'Ring', ['SignalID'], 'EventParam')
273     df = copy_down(df, 31, 'CycleStart', ['SignalID'], 'Timestamp')
274 else:
275     df['Ring'] = None
276     df['CycleStart'] = None
277 if 131 in codes:
278     df = copy_down(df, 131, 'CoordPattern', ['SignalID'], 'EventParam')
279 else:
280     df['CoordPattern'] = None
281 if 132 in codes:
282     df = copy_down(df, 132, 'CycleLength', ['SignalID'], 'EventParam')
283 else:
284     df['CycleLength'] = None
285 if 316 in codes:
286     df = copy_down(df, 316, 'ActualCycleLength', ['SignalID'], 'EventParam') # New 7/20/21
287 else:
288     df['ActualCycleLength'] = None
289 if 317 in codes:
290     df = copy_down(df, 317, 'ActualNaturalCycleLength', ['SignalID'], 'EventParam') # New 7/20/21
291 else:
292     df['ActualNaturalCycleLength'] = None
293 if 133 in codes:
294     df = copy_down(df, 133, 'CycleOffset', ['SignalID'], 'EventParam')
295 else:
296     df['CycleOffset'] = None
297 if 318 in codes:
298     df = copy_down(df, 318, 'ActualCycleOffset', ['SignalID'], 'EventParam') # New 7/20/21
299 else:
300     df['ActualCycleOffset'] = None
301 if 150 in codes:
302     df = copy_down(df, 150, 'CoordState', ['SignalID'], 'EventParam')
303 else:
304     df['CoordState'] = None
305 if 173 in codes:
306     df = copy_down(df, 173, 'FlashStatus', ['SignalID'], 'EventParam')
307 else:
308     df['FlashStatus'] = None
309
310 split_eventcodes = list(range(134,150))
311 df = create_new_grouping_field(df, split_eventcodes, 'EventCode', 'Phase', lambda x: x-133)
312
313 actual_split_eventcodes = list(range(300, 316))
314 df = create_new_grouping_field(df, actual_split_eventcodes, 'EventCode', 'Phase', lambda x: x-299)
315
316 ped_wait_eventcodes = list(range(612, 652))
317 df = create_new_grouping_field(df, ped_wait_eventcodes, 'EventCode', 'Phase', lambda x: x-611)

```

```

318
319 phase_eventcodes = list(range(0,25)) + list(range(41,50)) + [151]
320 df = create_new_grouping_field(df, phase_eventcodes, 'EventParam', 'Phase')
321 df = copy_down(df, 0, 'PhaseStart', ['SignalID','Phase'], 'Timestamp')
322 df = copy_down(df, [1,8,10], 'Interval', ['SignalID','Phase'], 'EventCode', apply_to_timestamp='group')
323 df.loc[df.EventCode==4, 'TermType'] = 4
324 df.loc[df.EventCode==5, 'TermType'] = 5
325 df.loc[df.EventCode==6, 'TermType'] = 6
326
327 # TODO0: See if we can get mapping between Vehicle Detector ID and Phase using (82, 81) and (43, 44).
328 #       Seems we can.
329 # TODO0: See if we can get mapping between Pedestrian Detector ID and Phase using (90), (45)
330
331 df = copy_down(df, split_eventcodes, 'ProgrammedSplit', ['SignalID','Phase'], 'EventParam',
332 apply_to_timestamp='group')
333 df = copy_up(df, actual_split_eventcodes, 'RecordedSplit', ['SignalID','Phase'], 'EventParam',
334 apply_to_timestamp='group')
335 df = copy_down(df, ped_wait_eventcodes, 'PedWait', ['SignalID','Phase'], 'EventParam', apply_to_timestamp=
336 'all')
337
338 df = copy_down(df, [183, 184], 'PowerFailure', ['SignalID'], 'EventParam', off_eventcode=182)
339
340 # Pair up detector on/offes under eventcode 82
341 df = copy_up(df, 81, 'DetectorOff', ['SignalID','Detector'], 'Timestamp', apply_to_timestamp=None)
342 df.loc[df.EventCode != 82, 'DetectorOff'] = np.nan
343 df['DetectorOff'] = pd.to_datetime(df['DetectorOff'])
344 df['DetectorDuration'] = (df['DetectorOff'] - df['Timestamp'])/pd.Timedelta(1, 's')
345
346 # Pair up ped input on/offes under eventcode 90
347 df = copy_up(df, 89, 'PedInputOff', ['SignalID','Detector'], 'Timestamp', apply_to_timestamp=None)
348 df.loc[df.EventCode != 90, 'PedInputOff'] = np.nan
349 df['PedInputOff'] = pd.to_datetime(df['PedInputOff'])
350 df['PedInputDuration'] = (df['PedInputOff'] - df['Timestamp'])/pd.Timedelta(1, 's')
351
352 df = df[~df.EventCode.isin([43, 44, 81, 89])]
353
354 # Pair up phase call on/offes under eventcode 43
355 #df = copy_up(df, 44, 'PhaseCallOff', ['SignalID','Detector'], 'Timestamp', apply_to_timestamp=None)
356 #df.loc[df.EventCode != 43, 'PhaseCallOff'] = np.nan
357 #df['PhaseCallDuration'] = (df['PhaseCallOff'] - df['Timestamp'])/pd.Timedelta(1, 's')
358 #df = df[~df.EventCode.isin([43, 44,81,89])]
359
360 # TODO0: Need a way to account for multiple detectors, inputs, etc. that overlap.
361 #       Add a new column for each? e.g., Detector1, Detector2, etc.?
362
363 # Possible update to copy_down for phase interval status. but needs to be grouped by phase
364 #df.loc[df.EventCode.isin(eventcodes), 'Interval'] = df.loc[df.EventCode.isin(eventcodes), 'EventCode']
365
366 df = df[df['Timestamp'].dt.date==date.date()].reset_index(drop=True)
367 df = df[['Timestamp','SignalID','EventCode','EventParam','date',
368 'Ring','CycleStart','CoordPattern','CoordState',
369 'CycleLength','ActualCycleLength','ActualNaturalCycleLength','CycleOffset','ActualCycleOffset',
370 'Phase','PhaseStart','Interval','TermType','ProgrammedSplit','RecordedSplit',
371 'Detector','DetectorFault','DetectorOff','DetectorDuration',
372 'PedInput','PedWait','PedInputOff','PedInputDuration']]
373
374 print('{} | {} done.'.format(date_str, s))
375
376 df.to_parquet(f's3://{events_bucket}/atspm_wide/date={date_str}/atspm_wide_{s}_{date_str}.parquet')
377
378 return df
379

```

```

380 def get_signalids(bucket, prefix):
381
382     s3 = boto3.client('s3')
383     paginator = s3.get_paginator('list_objects')
384
385     # Create a PageIterator from the Paginator
386     page_iterator = paginator.paginate(
387         Bucket=bucket,
388         Prefix=prefix)
389
390     for contents in [page['Contents'] for page in page_iterator]:
391         keys = [content['Key'] for content in contents]
392         for key in keys:
393             try:
394                 signalid = re.search('atspm_(.+?)_', key).group(1)
395             except AttributeError:
396                 signalid = ''
397
398             yield signalid
399
400
401
402 def get_channel_phase_mapping(df, channel_num=82, phase_num=43):
403
404     df = df.rename(columns={'TimeStamp': 'Timestamp'})
405     df = df.sort_values(['SignalID', 'Timestamp', 'EventCode', 'EventParam']).reset_index(drop=True)
406
407     is_phase = df.EventCode==phase_num
408     is_channel = df.EventCode==channel_num
409     df.loc[is_phase, 'Timestamp'] = df.loc[is_phase, 'Timestamp'] - timedelta(seconds=0.1)
410
411     dfc = df[is_channel].set_index(['SignalID', 'Timestamp'])['EventParam']
412     dfp = df[is_phase].set_index(['SignalID', 'Timestamp'])['EventParam']
413
414     dfcp = (pd.merge(left=dfc, right=dfp, on=['SignalID', 'Timestamp'], how='outer', suffixes=['_c', '_p'])
415             .dropna('EventParam_p')
416             .apply(tuple, axis=1))
417
418     dfcp = (dfcp[dfcp.groupby(level=['SignalID', 'Timestamp']).transform('count') == 1]
419             .reset_index(level='Timestamp', drop=True)
420             .drop_duplicates()
421             .sort_values())
422
423     dfcp = (pd.DataFrame.from_records(
424         dfcp,
425         columns=['Detector', 'CallPhase'],
426         index=dfcp.index)
427             .astype(int)
428             .reset_index(drop=False))
429
430     return dfcp
431
432
433
434
435
436
437
438
439
440
441
442
443
444

```

```

445 if __name__ == '__main__':
446
447     if len(sys.argv) > 1:
448         start_date = sys.argv[1]
449         end_date = sys.argv[2]
450     else:
451         #start_date = '2021-07-14'
452         #end_date = '2021-07-14'
453         sys.exit('Need start_date and end_date as command line parameters')
454
455     if start_date == 'yesterday':
456         start_date = (datetime.today() - timedelta(days=1)).strftime('%Y-%m-%d')
457     if end_date == 'yesterday':
458         end_date = (datetime.today() - timedelta(days=1)).strftime('%Y-%m-%d')
459
460
461     dates = pd.date_range(start_date, end_date, freq='1D')
462
463     for date_ in dates:
464         t0 = time.time()
465
466         date0_ = date_ - pd.Timedelta(1, unit='D')
467         date_str = date_.strftime('%Y-%m-%d')
468         print(date_str)
469
470         signalids = get_signalids(events_bucket, prefix=f'atspm/date={date_str}')
471         det_config = get_det_configs(config_bucket, [date0_, date_])
472
473         # with get_context('spawn').Pool(processes=os.cpu_count()-1) as pool:
474         with Pool(os.cpu_count()-1) as pool:
475             pool.starmap_async(widen, itertools.product(signalids, [date_], [det_config], ['s3']))
476             pool.close()
477             pool.join()
478
479         print(f'{len(signalids)} signals in {round(time.time()-t0, 1)} seconds.')
480
481 #
482 df[~df.DetectorDuration.isna()].groupby(['SignalID', 'CycleStart', 'Phase', 'Interval']).count()['Timestamp'].unstack('Interval', fill_value=0)
483 # df[~df.TermType.isna()][['SignalID', 'CycleStart', 'Phase', 'TermType']]

```