

D:\NHRI\AMR\amr-main\explore\pathogen_paper\get_burden_estimates_draws.py

```
1
2 import re
3 import os
4 import getpass
5 import glob
6 import argparse
7 import pandas as pd
8 import numpy as np
9 from pathlib import Path
10 import importlib
11 import ipdb
12 from datetime import date
13 from get_draws.api import get_draws
14
15 from db_queries import get_cause_metadata
16 from cod_prep.downloaders import get_pop, get_current_location_hierarchy
17 from cod_prep.claude.configurator import Configurator
18 import db_queries.api.public as db
19
20 from cod_prep.utils import print_log_message
21 from cod_prep.downloaders import get_ages
22 from amr_prep.utils.amr_io import get_amr_results
23 from mcod_prep.utils.mcause_io import get_mcause_results
24 from db_queries import get_outputs as go
25 from db_queries import get_population
26
27 # directories
28
29 DIR = 'FILEPATH'
30 repo_dir = 'FILEPATH'
31 MAP_DIR = 'FILEPATH'
32
33 DRAWS_DIR = 'FILEPATH'
34
35 CONF = Configurator('standard')
36 LSV_ID = CONF.get_id('location_set_version')
37 lh = get_current_location_hierarchy(location_set_version_id=LSV_ID)
38 ch = get_cause_metadata(cause_set_id=3, gbd_round_id=6, decomp_step="step4")
39
40 locs = lh.loc[lh["level"].isin([1,3]), "location_id"].unique().tolist() + [1]
41 ages = list(range(1,21))+[30, 31, 32, 235, 22, 27]
42 measures = [1,4]
43 metrics = [1,2]
44 sexes = [3]
45
46
47 gbd_unadj_causes_ids =[1027, 349, 843, 398, 405, 354, 355, 356, 397, 362, 363, 364, 936,
48 408]
49 gbd_adjusted_fpath = f'{DIR}gbd_causes_adjusted.csv'
50 gbd_adjusted_draws = f'{DIR}gbd_causes_adjusted_draws.csv'
51
52 pops = get_pop(release_id=6, pop_run_id=192)
```

```

52 pops = pops[pops.location_id.isin(locs)]
53
54 stomach_cancer_path= (f'{DIR}/stomach_cancer_adj.csv')
55 stomach_cancer_draws = (f'{DIR}/stomach_cancer_adj_draws.csv')
56
57 def pull_amr_causes_draws_u5():
58
59     df_ftl_u5_raw = get_amr_results(
60         process="summarize_burden",
61         burden_type='fatal',
62         draws=True,
63         year_id=2019,
64         cause_id=294,
65         age_group_id=[2,3,4,5],
66         location_id = locs,
67         measure_id=[1],
68         metric_id=[1,2], #number
69         sex_id=3,
70         hosp="all"
71     )
72
73     df_nonftl_u5_raw = get_amr_results(
74         process="summarize_burden",
75         burden_type='nonfatal',
76         draws=True,
77         year_id=2019,
78         cause_id=294,
79         age_group_id=[2,3,4,5],
80         location_id = locs,
81         measure_id=[2],
82         metric_id=[1,2], #number
83         sex_id=3,
84         hosp="all"
85     )
86
87     return pd.concat([df_ftl_u5_raw, df_nonftl_u5_raw])
88
89 def pull_amr_causes_draws_all_ages():
90
91     df_ftl_aa_raw = get_amr_results(
92         process="summarize_burden",
93         burden_type='fatal',
94         draws=True,
95         year_id=2019,
96         cause_id=294,
97         age_group_id=[22],
98         location_id = locs,
99         measure_id=[1],
100        metric_id=[1,2], #number
101        sex_id=3,
102        hosp="all"
103    )
104
105    df_nonftl_aa_raw = get_amr_results(

```

```

106     process="summarize_burden",
107     burden_type='nonfatal',
108     draws=True,
109     year_id=2019,
110     cause_id=294,
111     age_group_id=[22],
112     location_id = locs,
113     measure_id=[2],
114     metric_id=[1,2],
115     sex_id=3,
116     hosp="all"
117 )
118
119     return pd.concat([df_ftl_aa_raw, df_nonftl_aa_raw])
120
121 def get_additional_gbd_causes_draws(gbd_causes=gbd_unadj_causes_ids):
122
123     '''
124     Get draws for all gbd unadjusted causes, calculate counts/rates for Deaths, calculate
125     count/rates DALYs
126
127     parameters:
128     gbd_cause - a list of gbd causes that don't need additional chain/UCoD adjustment
129
130     returns:
131     a dataframe of 1000 draws, for Deaths and DALYs, counts and rates, all ages and under 5
132     age groups
133     '''
134
135     def check_metric_ids(df):
136         '''
137         check the both rates and counts metrics are present
138         '''
139         assert 1 in df.metric_id.unique(), 'no counts are present'
140         assert 3 in df.metric_id.unique(), 'no rates are present'
141
142     # get dalys rates
143
144     cc = get_draws("cause_id", gbd_causes,
145                   source="codcorrect",
146                   metric_id=[1],
147                   measure_id=[1,4],
148                   release_id=6,
149                   location_id=locs,
150                   sex_id=3,
151                   year_id=2019,
152                   age_group_id=[1,22])
153
154     dn = get_draws("cause_id", gbd_causes,
155                   source="dalynator",
156                   metric_id=[1],
157                   release_id=6,
158                   location_id=locs,

```

```

158         sex_id=3,
159         year_id=2019,
160         age_group_id=[1,22])
161
162
163
164     pops = get_population(age_group_id=[1, 22], location_id=locs, year_id=2019, sex_id=3,
165 release_id=6)
166     draw_cols = [f'draw_{i}' for i in range(1000)]
167
168     # Get DALYs rates
169     dn = dn.merge(pops[["age_group_id", "location_id", "population"]], how = 'left', on=
170 ["age_group_id", "location_id"], validate='m:1')
171     dn_rate = dn.copy()
172     for col in draw_cols:
173         dn_rate[col] = dn_rate[col]/dn_rate['population']
174     dn_rate['metric_id'] = 3
175     dalys = pd.concat([dn, dn_rate])
176     dalys = dalys.drop(['population', 'version_id'], axis=1)
177
178     assert len(dn.cause_id.unique())==14, 'not the correct number of gbd un-adjusted causes'
179     assert 1 in dn.age_group_id.unique()
180     assert 22 in dn.age_group_id.unique()
181
182     # Get Deaths rates
183     cc = cc.merge(pops[["age_group_id", "location_id", "population"]], how = 'left', on=
184 ["age_group_id", "location_id"], validate='m:1')
185     cc_rate = cc.copy()
186     for col in draw_cols:
187         cc_rate[col] = cc_rate[col]/cc_rate['population']
188     cc_rate['metric_id'] = 3
189     cc = pd.concat([cc, cc_rate])
190     cc = cc.drop(['population', 'version_id'], axis=1)
191     deaths = cc[cc.measure_id==1]
192
193     assert 1 in deaths.age_group_id.unique()
194     assert 22 in deaths.age_group_id.unique()
195
196     check_metric_ids(deaths)
197     check_metric_ids(dalys)
198
199     return pd.concat([deaths, dalys])
200
201
202 def add_viral_meningitis(df):
203     '''
204     Add AMR estimates for viral meningitis
205     '''
206     v_m = df.loc[(df.infectious_syndrome=='cns_infectious')&(df.pathogen=='virus')]
207     v_m['pathogen'] = 'viral_meningitis'
208
209     # filter "all" infectious syndromes

```

```

209     df = df.loc[df.infectious_syndrome=='all']
210
211     # append the viral meningitis
212     df = pd.concat([df,v_m])
213     # remove not needed columns
214     df.drop(columns=['infectious_syndrome', 'abx_class', 'abx_set', 'counterfactual',
'cause_id', 'hosp'], inplace=True)
215     return df
216
217 def add_pathogen_name_amr(df):
218
219     '''update pathogen name for the AMR causes based on the pathogen_metadata. It is not the
final step in assigning pathogen names for the results
220     '''
221     pathogen_mdata = pd.read_csv(f'FILEPATH')
222
223     df = df.merge(pathogen_mdata[['pathogen', 'pathogen_name_long']], how= 'left')
224     df.rename(columns = {'pathogen_name_long': 'pathogen_name'}, inplace=True)
225
226     df.loc[df.pathogen=='viral_meningitis', 'pathogen_name'] = 'Viral meningitis'
227     df['pathogen'] = df.pathogen_name
228     df = df.drop(['pathogen_name'], axis=1)
229
230     return df
231
232 def add_rates(df, pops_df=pops):
233
234     df = df.merge(pops_df, how='left', on =
['age_group_id', 'location_id', 'year_id', 'sex_id'], validate='m:1')
235
236     draw_cols = [f'draw_{i}' for i in range(1000)]
237     df_rate = df.copy()
238
239     for col in draw_cols:
240         df_rate[col] = df_rate[col]/df_rate['population']
241     df_rate['metric_id'] = 3
242
243     df =pd.concat([df, df_rate])
244     df = df.drop('population',axis=1)
245
246     return df
247
248 def format_pathogens_in_final_results(df):
249
250     '''
251     Format pathogen names to meet the requirements for the results
252
253     '''
254
255     ptgns = pd.read_csv(f'FILEPATH')
256
257     ptgns = ptgns[['pathogen_name', 'pathogen_name_updated']]
258     ptgn_list = list(ptgns.pathogen_name_updated.unique())
259     df = df.merge(ptgns, how='left', left_on='pathogen', right_on='pathogen_name')

```

```

260
261     df.pathogen = df['pathogen_name_updated']
262     df = df.sort_values(by='pathogen')
263
264     return df
265
266 def format_results(df, places = 1h[['location_id', 'location_name']]):
267
268     ages = get_ages()
269     measure_dict = {1: 'Deaths', 2: 'DALYs', 3: 'YLDs', 4: 'YLLs'}
270     metric_dict = {1: 'Count', 3: 'Rate'}
271
272     df['measure'] = df.measure_id.map(measure_dict)
273     df['metric'] = df.metric_id.map(metric_dict)
274     df = df.merge(ages[['age_group_id', 'age_group_name']], how='left', on='age_group_id',
275 validate='m:1')
276     df = df.merge(places, how='left', on='location_id', validate='m:1')
277
278     return df
279
280 def output_global_results(df, age_id, measure, metric):
281
282     cols = ['location_name', 'age_group_id', 'age_group_name', 'pathogen', 'mean_val',
283 'upper_val', 'lower_val', 'source', 'measure', 'metric']
284     df = df.loc[(df.measure==measure)&(df.metric==metric)&(df.age_group_id==age_id), cols]
285     df = df.sort_values(by='pathogen')
286     age_name = df.age_group_name.unique()[0]
287
288     return df
289
290 def save_global_results(df, dir_to_save, age_g):
291
292     assert age_g in ['Under_5', 'All_Ages'], "use 'Under_5', 'All_Ages'"
293     age_dict = {'Under_5': 1, 'All_Ages': 22}
294     age_id = age_dict[age_g]
295     today = date.today()
296     file_suffix = today.strftime("%m%d%y")
297
298     global_df = df.loc[df.location_name=='Global']
299     dalys_count = output_global_results(global_df, age_id, 'DALYs', 'Count')
300     deaths_count = output_global_results(global_df, age_id, 'Deaths', 'Count')
301     dalys_rate = output_global_results(global_df, age_id, 'DALYs', 'Rate_per_100_000')
302     deaths_rate = output_global_results(global_df, age_id, 'Deaths', 'Rate_per_100_000')
303
304     with pd.ExcelWriter(f'{dir_to_save}global_{age_g}_burden_estimates{file_suffix}.xlsx')
305 as writer:
306
307         # use to_excel function and specify the sheet_name and index
308         # to store the dataframe in specified sheet
309         dalys_count.to_excel(writer, sheet_name=f"DALYs_Count_{age_g}", index=False,
310 float_format='%10.2f')

```

```

309     deaths_count.to_excel(writer, sheet_name=f"Deaths_Count_{age_g}", index=False,
float_format='%10.2f')
310     dalsys_rate.to_excel(writer, sheet_name=f"DALYs_Rate_{age_g}", index=False,
float_format='%10.2f')
311     deaths_rate.to_excel(writer, sheet_name=f"Deaths_Rate_{age_g}", index=False,
float_format='%10.2f')
312
313
314 def output_regional_results(df, age_id, measure, metric):
315
316     cols = ['location_name', 'age_group_id', 'age_group_name', 'pathogen', 'mean_val',
'upper_val', 'lower_val', 'source', 'measure', 'metric']
317     df = df.loc[(df.measure==measure)&(df.metric==metric)&(df.age_group_id==age_id),cols]
318
319
320     df.rename(columns={'upper_val':f'{measure} {metric} upper',
321                       'lower_val': f'{measure} {metric} lower'},inplace=True)
322     df = df.sort_values(by='pathogen')
323     df = df.pivot(index = 'pathogen', columns='location_name', values=[ 'mean_val',
f'{measure} {metric} lower', f'{measure} {metric} upper'])
324
325     return df
326
327 def save_regional_results(df, dir_to_save, age_g):
328
329     '''
330     output all the results by regions
331     argunemts:
332         df: results dataframe
333         dir_to_save: the output directory
334         age_g: age group
335
336     '''
337
338     today = date.today()
339     file_suffix = today.strftime("%m%d%y")
340
341     assert age_g in ['Under_5', 'All_Ages'], "use 'Under_5', 'All_Ages'"
342     age_dict = {'Under_5':1 , 'All_Ages':22}
343     age_id = age_dict[age_g]
344
345     region_ids = lh.loc[lh.level==1, 'location_id'].unique()
346     region_results = df.loc[df.location_id.isin(region_ids)]
347
348     df1 = output_regional_results(region_results, age_id, 'DALYs', 'Count')
349     df2 = output_regional_results(region_results, age_id, 'Deaths', 'Count')
350     df3 = output_regional_results(region_results, age_id, 'DALYs', 'Rate_per_100_000')
351     df4 = output_regional_results(region_results, age_id, 'Deaths', 'Rate_per_100_000')
352
353
354     with
pd.ExcelWriter(f'{dir_to_save}regional_{age_g}_burden_estimates_{file_suffix}.xlsx') as
writer:
355

```

```

356 # use to_excel function and specify the sheet_name and index
357 # to store the dataframe in specified sheet
358 df1.to_excel(writer, sheet_name=f"DALYs_Count_{age_g}", float_format='%10.3f')
359 df2.to_excel(writer, sheet_name=f"Deaths_Count_{age_g}", float_format='%10.3f')
360 df3.to_excel(writer, sheet_name=f"DALYs_Rate_{age_g}", float_format='%10.3f')
361 df4.to_excel(writer, sheet_name=f"Deaths_Rate_{age_g}", float_format='%10.3f')
362
363
364
365 if __name__ == '__main__':
366
367     gbd = get_additional_gbd_causes_draws()
368
369     gbd_adj = pd.read_csv(gbd_adjusted_draws)
370
371     gbd_adj['age_group_id'] = gbd_adj.age_group_name.map({'All Ages':22, 'Under 5':1 })
372
373     amr_aa = pull_amr_causes_draws_all_ages()
374     amr_u5 = pull_amr_causes_draws_u5()
375
376     sc = pd.read_csv(stomach_cancer_draws)
377
378     # get the list of amr pathogens we estimate
379     all_amr_pathogens = amr_aa.loc[amr_aa.infectious_syndrome=='all', 'pathogen'].unique()
380
381     exclude_pathogens =
382     ['(none_estimated)', 'neisseria_gonorrhoeae', 'all', 'virus', 'other', 'mycobacterium_tuberculosis']
383     amr_pathogens = [i for i in all_amr_pathogens if i not in exclude_pathogens] +
384     ['viral_meningitis']
385     amr_pathogens.sort()
386
387     amr_aa = add_viral_meningitis(amr_aa)
388     amr_u5 = add_viral_meningitis(amr_u5)
389
390     # filter by pathogens
391
392     amr_aa = amr_aa.loc[amr_aa.pathogen.isin(amr_pathogens)]
393     amr_u5 = amr_u5.loc[amr_u5.pathogen.isin(amr_pathogens)]
394
395     # aggregate under 5
396
397     amr_u5_agg = amr_u5.groupby([ 'location_id', 'pathogen', 'sex_id',
398     'measure_id', 'year_id', 'metric_id']).sum().reset_index()
399     amr_u5_agg['age_group_id'] = 1
400
401     amr = pd.concat([amr_aa, amr_u5_agg])
402     amr = add_pathogen_name_amr(amr)
403
404     gbd['source'] = 'gbd'
405     amr['source'] = 'amr'
406     gbd_adj['source'] = 'gbd_adjusted'
407
408     gbd_adj['metric_id'] = 1

```



```

406     gbd_adj['year_id'] = 2019
407     gbd_adj['sex_id'] = 3
408
409     # keep only DALYs and Deaths
410     gbd_adj = gbd_adj[gbd_adj.measure_id.isin([1,2])]
411
412     sc['source'] = 'stomach_cancer'
413     sc['year_id'] = 2019
414
415     gbd = gbd.merge(ch[['cause_id', 'cause_name']], how='left', on='cause_id',
validate='m:1')
416     gbd = gbd.rename(columns={'cause_name': 'pathogen'})
417
418     measure_dict = {1: 'Deaths',
419                     2: 'DALYs',
420                     3: 'YLDs',
421                     4: 'YLLs'}
422
423     metric_dict = {1: 'Count',
424                   3: 'Rate'}
425
426
427     # add rates
428     gbd_adj = add_rates(gbd_adj)
429     amr = add_rates(amr)
430     sc = add_rates(sc)
431
432     results = pd.concat([gbd, gbd_adj, amr, sc])
433     results = results.drop(['cause_id', 'age_group_name'], axis=1)
434
435     # update the pathogen names based on the needs of the paper
436     results = format_pathogens_in_final_results(results)
437     results = results.drop(['pathogen', 'pathogen_name'], axis=1)
438     results = results.rename(columns={'pathogen_name_updated': 'pathogen'})
439
440     draw_cols = [f'draw_{i}' for i in range(0,1000)]
441     results['mean_val'] = results[draw_cols].apply(np.mean, axis=1)
442     results['lower_val'] = results[draw_cols].apply(np.quantile,q=0.025, axis=1)
443     results['upper_val'] = results[draw_cols].apply(np.quantile,q=0.975, axis=1)
444
445     cols_to_keep = ['metric_id', 'measure_id', 'age_group_id', 'location_id', 'sex_id',
'year_id', 'source', 'pathogen', 'mean_val', 'lower_val', 'upper_val']
446
447     locs = lh.loc[lh["level"].isin([1]), "location_id"].unique().tolist()+[1]
448
449     results_rg = results.loc[results.location_id.isin(locs)]
450
451     cols_to_agg = ['metric_id', 'measure_id', 'age_group_id', 'location_id', 'sex_id',
'year_id']
452     results_rg = results_rg.groupby(cols_to_agg).sum().reset_index(0)
453
454
455     results_rg['mean_val'] = results_rg[draw_cols].apply(np.mean, axis=1)
456     results_rg['lower_val'] = results_rg[draw_cols].apply(np.quantile,q=0.025, axis=1)

```

```
457 results_rg['upper_val'] = results_rg[draw_cols].apply(np.quantile,q=0.975, axis=1)
458
459 totals = results_rg.reset_index()
460 totals['pathogen'] = 'total'
461
462 results = pd.concat([results, totals])
463
464 # remove draw_cols
465 results = results.drop(draw_cols, axis=1)
466
467 # add measure, metric, age_group_name, location_name (def format_results)
468 results = format_results(results)
469
470 # add rates per 100_000
471 results_rate = results[results.metric_id==3]
472
473 for col in ['mean_val', 'lower_val', 'upper_val']:
474     results_rate[col] = results_rate[col]*100_000
475
476 results_rate['metric'] = 'Rate_per_100_000'
477 results_rate['metric_id'] = np.nan
478 results = pd.concat([results, results_rate])
479
480 # save all results in default format
481 results.to_csv(f'{DIR}FILEPATH', index=False)
482
483 output_dir = DIR
484 age_string = 'All_Ages'
485 save_global_results(results, output_dir, age_string)
486
487 age_string = 'Under_5'
488 save_global_results(results, output_dir, age_string)
489
490 output_dir = DIR
491 age_string = 'All_Ages'
492 save_regional_results(results, output_dir, age_string)
493
494 age_string = 'Under_5'
495 save_regional_results(results, output_dir, age_string)
```