

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

```
1  """
2  Purpose: Adjust several GBD causes according to mcod scalars (based on causes in chain and
3  as underlying cause of death) and add uncertainty
4  """
5
6  Getting adjustment values for several GBD causes
7
8  import pandas as pd
9  import numpy as np
10 import dask.dataframe as dd
11 import ipdb
12
13 from db_queries import get_cause_metadata
14 from cod_prep.downloaders import get_current_cause_hierarchy, add_location_metadata,
15 add_cause_metadata, get_current_location_hierarchy, prep_child_to_available_parent_map,
16 get_all_child_cause_ids, add_age_metadata
17 from mcod_prep.utils.mcause_io import *
18 from mcod_prep.utils.nids import add_nid_metadata
19 from amr_prep.utils.amr_io import *
20 from db_tools import ezfuncs, query_tools
21 from db_queries import get_outputs as go
22 from get_draws.api import get_draws
23 from db_queries import get_population
24
25 from functools import reduce
26 import argparse
27
28 CONF = Configurator()
29
30 lh = get_current_location_hierarchy()
31 ch = get_cause_metadata(cause_set_id=3, gbd_round_id=6, decomp_step="step4")
32
33 DIR = 'FILEPATH'
34 mcod_final = pd.read_csv(f'{DIR}FILEPATH')
35
36 all_u5 = [2,3,42,1,388,34,238,389,5,4]
37 gbd2019_u5 = [2,3,4,5]
38
39 # group ages for over 5
40 o5 = list(range(6,21)) + [30,31,32,235]
41
42 def change_names_of_pathogens(df):
43     """
44     change names of cause based on cause id
45     """
46     df.loc[df.cause_id==1026, 'cause_name'] = 'Total burden related to hepatitis B'
47     df.loc[df.cause_id==959, 'cause_name'] = 'Invasive Non-typhoidal Salmonella (iNTS)'
48     df.loc[df.cause_id==396, 'cause_name'] = 'Gonococcal infection'
49     df.loc[df.cause_id==320, 'cause_name'] = "Paratyphoid fever"
50     return df
51
52 def fix_ages(df, draws=True, keep_all=False):
```

```

50     '''
51     create under 5 and over 5 age categories
52     '''
53
54     keep_col1 = ["agg_age_group_id", "measure_id", "cause_id", "location_id"]
55     keep_col2 = ["age_group_id", "measure_id", "cause_id", "location_id"]
56
57     if draws == True:
58         keep_col1 += ["draws"]
59         keep_col2 += ["draws"]
60
61     df.loc[df["age_group_id"].isin(all_u5), "agg_age_group_id"] = 1
62     df.loc[df["age_group_id"].isin(o5), "agg_age_group_id"] = 192
63
64     df_u5 = df.loc[df["agg_age_group_id"] == 1]
65     df_o5 = df.loc[df["agg_age_group_id"] == 192]
66
67     if keep_all == True:
68         df_all = df.loc[df["age_group_id"] == 22]
69
70     df_u5 = df_u5.groupby(keep_col1, as_index=False)["val"].sum()
71     df_o5 = df_o5.groupby(keep_col1, as_index=False)["val"].sum()
72
73     df_u5.rename(columns={"agg_age_group_id": "age_group_id"}, inplace=True)
74     df_o5.rename(columns={"agg_age_group_id": "age_group_id"}, inplace=True)
75
76     if keep_all == True:
77         df_fixed = pd.concat([df_all, df_u5, df_o5])
78     else:
79         df_fixed = pd.concat([df_u5, df_o5])
80
81
82     df_fixed = add_age_metadata(df_fixed, "age_group_name")
83
84     return df_fixed
85
86 def output_scalars(tot_draws_df):
87     '''
88     save mean values and UI for the scalars used for gbd adjustment
89     '''
90
91
92     scalars = tot_draws_df[['draws', 'scalar', 'age_group_name', 'pathogen']]
93
94     scalars['mean_val'] = scalars.groupby(by= ['age_group_name', 'pathogen'])
95     ['scalar'].transform(np.mean)
96     scalars['lower_val'] = scalars.groupby(by= ['age_group_name', 'pathogen'])
97     ['scalar'].transform(np.quantile, q=0.025)
98     scalars['upper_val'] = scalars.groupby(by= ['age_group_name', 'pathogen'])
99     ['scalar'].transform(np.quantile, q=0.975)
100     scalars =
101     scalars.drop_duplicates(['age_group_name', 'pathogen', 'mean_val', 'lower_val', 'upper_val'])
102     scalars.drop(['draws', 'scalar'], axis=1, inplace=True)

```

```

100     scalars['scalar'] = scalars.apply(lambda x: f"{str(round(x.mean_val,2))}
({str(round(x.lower_val,2))} - {str(round(x.upper_val,2))})", axis=1)
101
102     scalars.rename(columns= {'age_group_name': 'Age', 'pathogen': 'Pathogen'},
inplace=True)
103     scalars.drop(['mean_val', 'lower_val', 'upper_val'], axis=1, inplace=True)
104
105     scalars.to_csv(f'{DIR}FILEPATH', index=False, float_format='%.2f')
106
107
108 def check_results(df, locs):
109
110     '''
111     check that all need columns are present in the results
112     '''
113     cols = ['location_id', 'pathogen', 'age_group_name', 'measure_id'] + draw_name
114     for col in cols:
115         assert col in df.columns, f'{col} is not present in the result dataframe'
116
117     # check it has all age groups
118     for i in ['All Ages', 'Under 5']:
119         assert i in df.age_group_name.unique(), f'{i} is missing from the dataframe'
120
121     #check the location number is correct
122     assert len(locs) == len(df.location_id.unique()), 'the number of locations is not
consistent'
123
124
125
126 if __name__ == '__main__':
127
128     parser = argparse.ArgumentParser()
129     parser.add_argument("--locs", help= "'global' to for global locations, 'region' for
regions, and 'country' for countries. if not specified, runs for regions and
global", type=str, default='all')
130     parser.add_argument("--scalars", help= '"save" to output scalars')
131     arg = parser.parse_args()
132
133
134     if arg.locs == 'global':
135         locs = [1]
136         loc_type = 'global'
137     elif arg.locs == 'region':
138         locs = lh.loc[lh["level"].isin([1]), "location_id"].unique().tolist()
139         loc_type = 'regions'
140     elif arg.locs == 'country':
141         locs = lh.loc[lh["level"]==3, "location_id"].unique().tolist()
142         loc_type = 'countries'
143     else:
144         locs = lh.loc[lh["level"].isin([2,1]), "location_id"].unique().tolist()+[1]
145         loc_type = 'regions_global'
146
147
148

```

```

149 #list of GBD pathogens of interest (POI) to calculate uncertainty for
150 poi = [
151     "Tuberculosis",
152     'African trypanosomiasis',
153     "HIV/AIDS",
154     "Malaria",
155     "Pertussis", #Whooping Cough
156     "Syphilis",
157     "Measles",
158     "Acute hepatitis A",
159     "Dengue",
160     "Tetanus",
161     "Hepatitis B",
162     "Varicella and herpes zoster",
163     "Rabies",
164     "Schistosomiasis",
165     "Chagas disease",
166     "Visceral leishmaniasis",
167     "Diphtheria",
168     "Yellow fever",
169     "Ascariasis",
170     "Acute hepatitis E",
171     "Cystic echinococcosis",
172     "Cysticercosis",
173     "Cervical cancer",
174     "Other neglected tropical diseases",
175     "Zika virus",
176     "Neisseria gonorrhoeae",
177 ]
178
179 mcod_poi = mcod_final.loc[mcod_final["pathogen"].isin(poi)]
180
181 # Calculate uncertainty
182 x = 1000
183 draw_name = [f'draw_{i}' for i in list(range(0, x))]
184
185 draw_dfs = []
186
187 np.random.seed(42)
188 for path in mcod_poi["pathogen"].unique().tolist():
189     for age in mcod_poi.loc[mcod_poi["pathogen"] == path,
190 "age_group_name"].unique().tolist():
191         path_df = mcod_poi.loc[(mcod_poi["pathogen"] == path) &
192 (mcod_poi["age_group_name"] == age)]
193
194         n = float(path_df["Deaths where UCoD"] + path_df["Deaths in Chain"])
195         if n != 0:
196             p = float(path_df["Deaths where UCoD"]/n)
197
198             draws = (1/(np.random.binomial(n, p, x)/n))
199
200             path_draws = pd.DataFrame(list(zip(draw_name, draws)),
201                                     columns=['draws', 'scalar'])
202             path_draws["age_group_name"] = age

```

```

201         path_draws["pathogen"] = path
202
203         draw_dfs.append(path_draws)
204
205     tot_draws = pd.concat(draw_dfs)
206     tot_draws.replace([np.inf, -np.inf], 1, inplace=True)
207
208     tot_draws = tot_draws.fillna(1)
209
210
211     mcod_poi["total"] = mcod_poi["Deaths where UCoD"] + mcod_poi["Deaths in Chain"]
212     mcod_poi_for_vet = mcod_poi.groupby(["pathogen", "age_group_name"], as_index=False)
213     ["total"].sum()
214
215     mcod_poi_for_vet = mcod_poi_for_vet.pivot(index=["pathogen"], columns="age_group_name",
216     values="total")
217     mpv = mcod_poi_for_vet.reset_index()
218
219     # causes that shouldn't have any correction (scalar = 1)
220     no_correction = []
221
222     # causes that should have a all ages scalar used for all age and for under 5
223     to_correct_u5 = []
224
225     # causes that should have a all age scalar for over 5
226     to_correct_o5 = []
227
228     cut_off = 200
229
230     for i in mpv.pathogen.unique():
231         if i not in ['Acute hepatitis E', 'Cervical cancer']:
232             if mpv.loc[mpv.pathogen==i, 'All Ages'].values[0]<200:
233                 no_correction.append(i)
234                 continue
235             elif mpv.loc[mpv.pathogen==i, 'Under 5'].values[0]<200:
236                 to_correct_u5.append(i)
237                 continue
238
239     # check if all_age corrections needed for over 5:
240     for i in mpv.pathogen.unique():
241         if i not in ['Acute hepatitis E', 'Cervical cancer']:
242             if (mpv.loc[mpv.pathogen==i, 'All Ages'].values[0]>=200)&
243             (mpv.loc[mpv.pathogen==i, '5 plus'].values[0]<200):
244                 to_correct_o5.append(i)
245
246
247     print(f"we are not correcting the following causes: {no_correction}")
248     print(f"we are using the all age scalar the following causes for under 5:
249     {to_correct_u5}")
250     print(f"we are using the all age scalar the following causes for over 5:
251     {to_correct_o5}")
252
253
254     # no correction for under 5 and all ages:
255     for i in no_correction:
256         tot_draws.loc[tot_draws["pathogen"] == i, "scalar"] = 1

```

```

250
251 # to correct using only all age scalars
252 for i in to_correct_u5:
253     tot_draws.loc[(tot_draws["pathogen"] == i)
254                   & (tot_draws["age_group_name"] == "Under 5"), "scalar"] =
tot_draws.loc[(tot_draws["pathogen"] == i)
255               & (tot_draws["age_group_name"] == "All Ages"), "scalar"]
256
257 for i in to_correct_o5:
258     tot_draws.loc[(tot_draws["pathogen"] == i)
259                   & (tot_draws["age_group_name"] == "5 plus"), "scalar"] =
tot_draws.loc[(tot_draws["pathogen"] == i)
260               & (tot_draws["age_group_name"] == "All Ages"), "scalar"]
261
262
263 # remove Cervical cancer for under 5
264 tot_draws = tot_draws.loc[~((tot_draws["pathogen"] == "Cervical cancer")
265                             & (tot_draws["age_group_name"] == "Under 5"))]
266
267 # Adjust for Hepatitis E
268 tot_draws.loc[(tot_draws["pathogen"] == "Acute hepatitis E")
269               & (tot_draws["age_group_name"] == "Under 5"), "scalar"] = 1
270 # add cause_id to tot_draws
271 tot_draws = tot_draws.merge(ch[['cause_name', 'cause_id']], how='left',
left_on='pathogen', right_on='cause_name')
272
273 # Name and cause id adjustments
274 tot_draws.loc[tot_draws.pathogen=='Neisseria gonorrhoeae', 'cause_id'] =
ch.loc[ch.cause_name=='Gonococcal infection', 'cause_id'].unique()[0]
275 tot_draws.loc[tot_draws.pathogen=='Hepatitis B', 'cause_id'] =
ch.loc[ch.acause=='total_hep_b_reporting', 'cause_id'].unique()[0]
276 tot_draws.loc[tot_draws.pathogen=='Pertussis', 'cause_id'] =
ch.loc[ch.acause=='whooping', 'cause_id'].unique()[0]
277
278 if arg.scalars == 'save':
279     output_scalars(tot_draws)
280
281 gbd_causes = tot_draws.cause_id.unique()
282
283 assert len(gbd_causes)==26
284
285 # get draws (counts) for Deaths and YLL (years of life lost)
286 cc = get_draws("cause_id", gbd_causes,
287               source="codcorrect",
288               metric_id=1,
289               measure_id=[1,4],
290               release_id=6,
291               location_id=locs,
292               sex_id=3,
293               year_id=2019,
294               age_group_id=gbd2019_u5 + o5)
295 cc = cc.fillna(0)
296
297 # Load in YLDs (rates)

```

```

298 como = get_draws("cause_id", gbd_causes,
299                   source="como",
300                   metric_id=3,
301                   measure_id=[3],
302                   release_id=6,
303                   location_id=locs,
304                   sex_id=3,
305                   year_id=2019,
306                   age_group_id=gbd2019_u5 + o5)
307
308
309 # Melt draws for gbd
310 IDs = ["measure_id",
311        "location_id",
312        "metric_id",
313        "age_group_id",
314        "cause_id"
315       ]
316
317 print("melting!")
318 cc = cc.melt(
319     id_vars=IDs,
320     value_vars=draw_name, var_name='draws', value_name='val'
321 )
322 cc["val"] = cc["val"].fillna(0)
323
324 como = como.melt(
325     id_vars=IDs,
326     value_vars=draw_name, var_name='draws', value_name='val'
327 )
328
329
330 pops = get_population(age_group_id=gbd2019_u5 + o5 + [22], location_id=locs,
331 year_id=2019, sex_id=3, release_id=6)
332
333 como = como.merge(pops[["age_group_id", "location_id", "population"]], how="left", on=
334 ["age_group_id", "location_id"], validate='m:1')
335
336 como_count = como.copy()
337 como_count["val"] = como_count["val"] * como_count["population"]
338 como_count["metric_id"] = 1
339 como_count["val"] = como_count["val"].fillna(0)
340
341 cc = fix_ages(cc, draws=True)
342 como_count = fix_ages(como_count, draws=True)
343
344 cc = add_cause_metadata(cc, "cause_name")
345 cc = change_names_of_pathogens(cc)
346
347 como_count = add_cause_metadata(como_count, "cause_name")
348 como_count = change_names_of_pathogens(como_count)
349
350 cc.rename(columns={"cause_name": "pathogen"}, inplace=True)
351 como_count.rename(columns={"cause_name": "pathogen", "val": "ylds"}, inplace=True)

```

```

350
351 # Fix for exceptions on GBD
352 cc = cc.loc[~((cc["pathogen"] == "Cervical cancer")
353             & (cc["age_group_name"] == "Under 5"))]
354
355 pathogen_dict = {'Non-typhoidal Salmonella': 'Invasive Non-typhoidal Salmonella (iNTS)',
356                 'Hepatitis B': 'Total burden related to hepatitis B',
357                 'Neisseria gonorrhoeae': 'Gonococcal infection'}
358
359 tot_draws.pathogen=tot_draws.pathogen.apply(lambda x: pathogen_dict[x] if x in
pathogen_dict.keys() else x)
360
361 to_fix = cc.merge(tot_draws, how="left", on=["draws", "cause_id",
"age_group_name"], suffixes= ['_cc', '_tot_draws'], validate='m:1')
362 to_fix = to_fix.drop(['pathogen_cc', 'cause_name'], axis=1)
363 to_fix = to_fix.rename(columns={'pathogen_tot_draws': 'pathogen'})
364
365 to_fix["fixed_fatal"] = to_fix["val"]*to_fix["scalar"]
366
367 ylls = to_fix.loc[to_fix["measure_id"] == 4]
368 ylls.rename(columns={"fixed_fatal": "fixed_ylls"}, inplace=True)
369
370 deaths = to_fix.loc[to_fix["measure_id"] == 1]
371 deaths.rename(columns={"fixed_fatal": "fixed_deaths"}, inplace=True)
372
373 id_cols = ["location_id", "pathogen", "draws", "age_group_name"]
374 measures = [como_count[id_cols + ["ylds"]],
375             ylls[id_cols + ["fixed_ylls"]],
376             deaths[id_cols + ["fixed_deaths", "scalar"]]]
377
378 final_fix = reduce(lambda left, right: pd.merge(left, right, on=id_cols,
379                                                how='outer'), measures)
380
381 final_fix["fixed_dalys"] = final_fix["fixed_ylls"] + final_fix["ylds"]
382
383
384 # create All Age category
385 final_fix_aa = final_fix.groupby(["location_id", "pathogen", "draws"], as_index=False)
["fixed_deaths",
386                                     "fixed_ylls",
387
"fixed_dalys"].apply(lambda x: x.sum())
388 final_fix_aa["age_group_name"] = "All Ages"
389
390 final_fix = pd.concat([final_fix, final_fix_aa])
391 final_fix = final_fix.fillna(0)
392
393 final_fix.drop("ylds", axis=1, inplace=True)
394 final_fix = pd.melt(final_fix, id_vars=["location_id", "pathogen", "draws",
"age_group_name", "scalar"],
395                     value_vars=["fixed_ylls", "fixed_deaths", "fixed_dalys"],
396                     var_name="measure_id",
397                     value_name='fixed_val',
398                     ignore_index=True)

```



```

399
400
401     measure_map = {
402         "fixed_ylls":4,
403         "fixed_deaths":1,
404         "fixed_dalys":2
405     }
406
407     final_fix["measure_id"] = final_fix["measure_id"].map(measure_map)
408
409     results= final_fix.drop('scalar', axis=1)
410
411     results = results.pivot(index= ['location_id',
412 'pathogen', 'age_group_name', 'measure_id'], columns='draws').reset_index()
413
414     draws_name=[f'draw_{i}' for i in range(1000)]
415
416     results.columns = results.columns.droplevel(1)
417
418     cols = []
419     count = 0
420     for column in results.columns:
421         if column == 'fixed_val':
422             cols.append(f'draw_{count}')
423             count+=1
424             continue
425         cols.append(column)
426
427     results.columns = cols
428     results = results[results.age_group_name.isin(['All Ages', 'Under 5'])]
429
430     assert len(results.pathogen.unique())==27)
431     check_results(results, locs)
432     results.to_csv(f'FILEPATH', index=False)

```