

D:\NHRI\AMR\amr-main\datasets\clean_SOURCE_1.py

```
1 import pandas as pd
2 import numpy as np
3 import argparse
4 import re
5 import sys
6 from cod_prep.claude.configurator import Configurator
7 from cod_prep.downloaders import get_ages, get_map_version, pretty_print
8 from cod_prep.utils import (report_if_merge_fail, print_log_message,
9                             clean_icd_codes)
10 from cod_prep.utils.formatting import ages
11 from mcod_prep.mcod_mapping import MCoDMapper
12 from amr_prep.utils.amr_io import save_amr_data
13
14 CONF = Configurator()
15 SOURCE = "SOURCE_NAME"
16 L_DIR = "FILEPATH"
17
18
19 def create_age_group_id(df):
20
21     df.loc[df.age < 0, "age"] = np.NaN
22
23     df['age_unit'] = 'year'
24     df.rename(columns={'Age': 'age'}, inplace=True)
25     age_formatter = ages.PointAgeFormatter()
26     df = age_formatter.run(df)
27
28     assert df['age_group_id'].notnull().values.all()
29
30     return df
31
32
33 def create_sex_id(df):
34     df['sex_id'] = df['sex'].apply(
35         lambda x: 1 if x == 'male' else (2 if x == 'female' else 9)
36     )
37     print(
38         f"{(len(df.loc[df['sex_id'] == 9]) / len(df)) * 100}% of rows are missing sex")
39     report_if_merge_fail(df, 'sex_id', 'sex')
40     return df
41
42
43 def mark_admissions_deaths(df):
44
45     discharge_map = pd.read_csv("FILEPATH")
46     discharge_dict = discharge_map.set_index("discharge_disp")[
47         "Died"].to_dict()
48
49     df["admissions"] = 1
50
51     df["deaths"] = df["discharge_disp"].map(discharge_dict)
```

```

52 df.loc[(df.deaths.isnull()) & (df.discharge_disp.str.contains(
53     "hospice", flags=re.IGNORECASE)), "deaths"] = 1
54 df.loc[df.discharge_disp.isin(["nan", "UNKNOWN", "-"]), "deaths"] = np.nan
55 return df
56
57
58 def apply_pathogen_severity_map(df, named_cols):
59
60     pathogen_severity = pd.read_excel("FILEPATH")
61     pathogen_dict = pathogen_severity.set_index("pathogen_name")["rank"].to_dict()
62
63     df[named_cols] = df[named_cols].apply(lambda x: x.map(
64         pathogen_dict)).fillna(100).astype(int)
65
66     df["severe_pathogen_int"] = df[named_cols].apply(lambda x: x.min(), axis=1)
67
68     reverse_dict = dict([value, key] for key, value in pathogen_dict.items())
69     df["severe_pathogen"] = df["severe_pathogen_int"].map(
70         reverse_dict)
71
72     df[named_cols] = df[named_cols].apply(lambda x: x.map(
73         reverse_dict))
74
75     df["severe_pathogen_int"] = df["severe_pathogen_int"].replace(100, np.NaN)
76
77     return df
78
79
80 def fix_multiple_deaths(df):
81     date_map = pd.read_csv("FILEPATH")
82     df = mark_admissions_deaths(df)
83
84     date_map.rename(columns={"date_id": "discharge_date_id",
85                             "date": "discharge_date"}, inplace=True)
86     df = df.merge(date_map, on="discharge_date_id", how="left")
87     df["discharge_date"] = pd.to_datetime(df.discharge_date)
88     report_if_merge_fail(df, "discharge_date_id", "date")
89     null = df.loc[df['patient_id'].isna(), ]
90     notnull = df.loc[~df['patient_id'].isna(), ]
91     notnull["max"] = notnull.groupby("patient_id", as_index=False)
92     ["discharge_date"].transform("max")
93     notnull["deaths"] = np.where(notnull["discharge_date"] == notnull["max"],
94     notnull["deaths"], 0)
95     df = pd.concat([null, notnull])
96
97     return df
98
99 def map_specimen_source(df):
100     specimen_detail = pd.read_excel("FILEPATH")
101     specimen_detail.rename(columns={"specimen_source_group": "specimen_source"},

```

```

102     df = df.merge(specimen_detail, on=["site_id", "specimen_id", "specimen_source"],
103 how="left")
104     df.loc[df.specimen_source_detail.isnull(), "specimen_source_detail"] =
105 df["specimen_source"]
106     report_if_merge_fail(df, "specimen_source_detail", "specimen_source")
107
108     df.rename(columns={'specimen_source_detail': 'raw_specimen'}, inplace=True)
109
110     return df
111
112 def apply_syndrome_severity_map(df):
113     priority_dict = {"peritonitis": 1,
114                      "lri": 2,
115                      "uti": 3,
116                      "sepsis/bacteremia": 4,
117                      "cellulitis": 5,
118                      "other": 6}
119
120     df['syndrome_num'] = df.groupby('admission_id')['specimen_syndrome'].transform(
121         lambda syndromes: [
122             f"specimen_syndrome_{i}" for i in range(len(syndromes))]
123     )
124     syndrome_df = df.pivot(
125         index='admission_id', columns='syndrome_num', values='specimen_syndrome'
126     ).reset_index()
127
128     syndrome_cols = [x for x in list(syndrome_df) if "specimen_syndrome" in x]
129     syndrome_df[syndrome_cols] = syndrome_df[syndrome_cols].apply(lambda x: x.map(
130         priority_dict)).fillna(100).astype(int)
131     syndrome_df = syndrome_df.assign(severe_syndrome_int=lambda x: x.min(axis=1))
132
133     reverse_dict = dict([value, key] for key, value in priority_dict.items())
134     syndrome_df["severe_syndrome"] = syndrome_df["severe_syndrome_int"].map(
135         reverse_dict)
136
137     df = df.merge(
138         syndrome_df[["admission_id", "severe_syndrome"]], on="admission_id", how="left")
139     drop_synds = df[df["specimen_syndrome"] != df["severe_syndrome"]].index
140     df.drop(drop_synds, inplace=True)
141
142     df.drop(columns="severe_syndrome", axis=1, inplace=True)
143
144     return df
145
146 def create_pathogen_cols(df):
147     df.rename(columns={'org_name': 'raw_pathogen'}, inplace=True)
148     df['raw_pathogen'] = df['raw_pathogen'].str.strip()
149
150     return df
151
152
153 def mark_estimated_pathogens(df, cols):

```

```

154 df["pathogen"] = np.where(df[cols].isnull().all(axis=1), "none", "other")
155 df["pathogen"] = np.where(df["severe_pathogen"].notnull(),
156                             df["severe_pathogen"], df["pathogen"])
157
158 return df
159
160
161 def map_cause_cols(df):
162     print_log_message("creating pathogens cols")
163     df = create_pathogen_cols(df)
164     diag_df = prep_diagnosis()
165
166     df = df.merge(diag_df, on="admission_id", how="left")
167     df.rename(columns={"primary_diagnosis_code": "cause"}, inplace=True)
168
169     print_log_message("mapping specimen source col")
170     df = map_specimen_source(df)
171
172     print_log_message("fixings patients with multiple 'deaths'")
173     df = fix_multiple_deaths(df)
174
175     for col in [x for x in list(df) if "cause" in x]:
176         df[col] = df[col].fillna("none")
177     df['cause'] = clean_icd_codes(df['cause'], remove_decimal=True)
178     df.loc[~df.cause.str.match(r"[A-Z]\d{2,5}"), 'cause'] = "none"
179
180     return df
181
182
183 def prep_diagnosis():
184     df = pd.read_csv("FILEPATH")
185     df['diagnosis_code'] = clean_icd_codes(df['diagnosis_code'], remove_decimal=True)
186     df.loc[~df.diagnosis_code.str.match(r"[A-Z]\d{2,5}"), 'diagnosis_code'] = "none"
187     df['cause_num'] = df.groupby('admission_id')['diagnosis_code'].transform(
188         lambda codes: [f"multiple_cause_{i}" for i in range(len(codes))]
189     )
190     df.loc[df['primary_diagnosis'], 'cause_num'] = 'cause'
191     df = df.pivot(
192         index='admission_id', columns='cause_num', values='diagnosis_code'
193     ).reset_index().fillna('none')
194     df.drop(columns="cause", inplace=True)
195     return df
196
197
198 def map_hosp(df):
199     date_map = pd.read_csv("FILEPATH")
200
201     admit_date_map = date_map.rename(columns={"date_id": "admit_date_id",
202                                                "date": "admit_date"})
203     collected_date_map = date_map.rename(columns={"date_id": "collected_date_id",
204                                                  "date": "collected_date"})
205
206     df = df.merge(admit_date_map, on='admit_date_id', how='left')
207     df = df.merge(collected_date_map, on='collected_date_id', how='left')

```

```

208 df['date_diff'] = (pd.to_datetime(df['collected_date']) -
209                    pd.to_datetime(df['admit_date'])).dt.days
210
211 df['hosp'] = 'hospital'
212 df.loc[df['date_diff'] <= 2, 'hosp'] = 'community'
213 df.loc[df['date_diff'].isna(), 'hosp'] = 'unknown'
214
215 return df
216
217
218 def clean_SOURCE():
219
220     df = pd.read_csv("FILEPATH")
221     cols = [col for col in df.columns if 'result_id' not in col]
222     df = df.drop_duplicates(cols)
223     df = create_age_group_id(df)
224     df = create_sex_id(df)
225     df = map_cause_cols(df)
226     df = map_hosp(df)
227     df.rename(columns={'drug_name': 'raw_antibiotic'}, inplace=True)
228     df['resistance'] = df.interpretation.map({
229         'Resistant': 'resistant',
230         'Susceptible': 'sensitive',
231         'Intermediate': 'resistant',
232         'Positive': 'resistant',
233         'Non Suceptible': 'resistant',
234         'Negative': 'sensitive'
235     }).fillna('unknown')
236
237     df["location_id"] = 102
238     df["nid"] = 410446
239     df["year_id"] = df["discharge_date"].apply(lambda x: x.year)
240     df['sample_id'] = (df['patient_id'].astype(str) + df['specimen_id'].astype(str))
241     df['cases'] = 1
242
243     df['hospital_type'] = 'unknown'
244     df['hospital_name'] = 'unknown'
245
246     cause_columns = [x for x in list(df) if "multiple_cause" in x] + ['cause']
247     demo_cols = ['nid', 'location_id', 'year_id', 'age_group_id', 'sex_id']
248     biology = ['raw_specimen', 'raw_pathogen', 'raw_antibiotic', 'resistance']
249     other_values = ['hosp', 'hospital_name', 'hospital_type', 'deaths', 'sample_id',
250 'cases']
251
252     lowercase_cols = ['raw_specimen', 'raw_pathogen', 'raw_antibiotic']
253
254     df = df[demo_cols + cause_columns + biology + other_values]
255
256     for col in lowercase_cols:
257         df[col] = df[col].str.lower().str.strip()
258     assert df.sample_id.notnull().values.all()
259
260     df = df.drop_duplicates()
261
262     return df

```

```
261 |
262 |
263 | if __name__ == "__main__":
264 |     df = clean_SOURCE()
265 |     save_amr_data(df, phase='unmapped', ldir=L_DIR)
266 |
```