

数据集清理

00a-prescriptions

读取和处理PIS OST处方数据，将变量重命名为统一的名称，并保存处理后的数据。

```
# Name of file - 05-PIS_OST_all_prescriptions.R ####
```

```
# Original author(s) - Morven Millar
```

```
# Original date - 15/12/2022
```

```
# Latest update author - Morven Millar
```

```
# Latest update date -
```

```
# Latest update description -
```

```
# Type of script - Data preparation of unaggregated file of PIS OST prescriptions
```

```
# Written/run on - R Studio Server
```

```
# Version of R - 3.6.1
```

```
# Description of content
```

```
# Script to call in the unaggregated file of all PIS OST prescriptions from Paid FY 2009/10 - 2021/22.
```

```
# 该脚本用于调用2009/10财年至2021/22财年所有PIS OST处方的未聚合文件。
```

```
# Approximate run time - <5mins
```

```
# 清除环境变量
```

```
rm(list = ls())
```

```
gc()
```

```
ok
```

```
ok# install packages
```

```
install.packages('tidyverse')
```

```
install.packages('dplyr')
```

```

# Load packages
library(tidyverse)
library(dplyr)
#设置工作目录,根据目录找到位置
setwd("")
getwd()
ost <- read_rds('OST_prescriptions.rds')
df01_OST_cohort <-
readRDS("/PHI_conf/SubstanceMisuse2/Topics/Prevalence/Opioids/Projects/2025-05-01-
FanZhang-MSc-OAT-cocaine/Analysis/PHS_raw_data/cohort_OST_prescriptions.rds")
#Converting Vars to match all files

ost <- ost %>%
  rename(PaidDate = paid_month_end_date,
         PrescrDate = prescribed_full_date,
         DispDate = dispensed_full_date,
         DOB = patient_date_of_birth,
         Gender = patient_sex,
         PatientsResidentsNHSBoard = patient_health_board_name,
         ItemCode = item_code)
range(ost$PaidDate)

write_rds(ost, "refined_data/OST_prescriptions.rds")

```

- 变量重命名:

- 使用 `dplyr` 包的 `rename` 函数, 将 `ost` 数据框中的变量重命名为统一的名称, 以便与所有文件匹配。
- 重命名的变量包括:
 - `PaidDate` (付账日期)
 - `PrescrDate` (处方日期)

- `DispDate` (发药日期)
 - `DOB` (患者出生日期)
 - `Gender` (患者性别)
 - `PatientsResidentsNHSBoard` (患者所在NHS委员会名称)
 - `ItemCode` (项目代码)
- 检查日期范围: 使用 `range(ost$PaidDate)` 检查 `PaidDate` 变量的范围。
 - 保存处理后的数据: 将处理后的 `ost` 数据保存为 `refined_data/OST_prescriptions.rds` 文件。

00b-strong_opioids_short

合并和处理与阿片类药物相关的数据集，标记强阿片类药物，并保存处理后的数据。

输入文件

1. 药物相关死亡 (DRDs) 数据集：
 - 文件名: `NRS_drug_related_deaths.rds`
 - 用途: 包含药物相关死亡的详细信息。代码中将其与额外的案例数据集合并。
2. 额外的药物相关死亡案例数据集：
 - 文件名: `NRS_drug_related_deaths_ADDITIONAL_CASES.rds`
 - 用途: 包含额外的药物相关死亡案例，用于扩展主数据集。
3. 强阿片类药物数据集：
 - 文件名: `lookup_strong_opioids.rds`
 - 用途: 包含强阿片类药物的相关信息，用于标记主数据集中的强阿片类药物案例。
4. 额外的强阿片类药物数据集：
 - 文件名: `lookup_strong_opioids_ADDITIONAL_CASES.rds`
 - 用途: 包含额外的强阿片类药物案例，用于扩展主强阿片类药物数据集。

输出文件

这段代码生成了以下输出文件：

1. 处理后的药物相关死亡数据集：

- 文件名： `refined_data/allDRDs.rds`
- 内容：经过变量转换、合并和标记后的药物相关死亡数据集。这个文件包含了主数据集和额外案例数据集的合并结果，以及强阿片类药物的标记。

Code block

```
1  # 清除环境变量 加载多个R包
2  rm(list=ls())
3  library(readr)
4  library(retroharmonize)
5  library(data.table)
6  library(tidyverse)
7  library(phsmethods)
8
9  # 设置工作目录
10 setwd(' ')
11 OpioidCohort <- read_rds("NRS_drug_related_deaths.rds")
12 ExtraOpioidCohort <- read_rds("NRS_drug_related_deaths_ADDITIONAL_CASES.rds")
13 # 检查NRSdrd_DOD (死亡日期) 的范围。
14 range(OpioidCohort$NRSdrd_DOD)
15 # 检查两个数据集的列名是否一致，确保合并。
16 colnames(OpioidCohort)==colnames(ExtraOpioidCohort)
17 OpioidCohort <- rbind(OpioidCohort,ExtraOpioidCohort)
18
19 # 检查合并
20 StrongOpioidCohort <- read_rds('lookup_strong_opioids.rds')
21 ExtraStrongOpioidCohort <-
22   read_rds('lookup_strong_opioids_ADDITIONAL_CASES.rds')
23 colnames(StrongOpioidCohort)==colnames(ExtraStrongOpioidCohort)
24 StrongOpioidCohort <- rbind(StrongOpioidCohort,ExtraStrongOpioidCohort)
25
26 # 将OpioidCohort数据集中的IAIN列转换为字符类型。
27 OpioidCohort <- OpioidCohort %>%
28   mutate(across(c("IAIN"), as.character))
29
30 # 将StrongOpioidCohort数据集中的IAIN列转换为字符类型，
31 # 并添加一个新列strong_opioid_flag，其值为1，用于标记强阿片类药物。
32 StrongOpioidCohort <- StrongOpioidCohort %>%
33   mutate(across(c("IAIN"), as.character)) %>%
34   mutate(strong_opioid_flag=1)
35
36
37 # 数据合并与标记
38 OpioidCohort <- OpioidCohort %>%
```

```

39         left_join(StrongOpioidCohort %>% select(IAIN,
strong_opioid_flag)) %>%
40
mutate(strong_opioid_flag=ifelse(is.na(strong_opioid_flag),0,strong_opioid_flag
))
41 # 统计标记数量
42 OpioidCohort %>% count(strong_opioid_flag)
43
44 getwd()
45 write_rds(OpioidCohort, "refined_data/allDRDs.rds")
46

```

00c-create_deaths_file

处理与药物相关死亡（DRDs）和全因死亡（ACM）的数据集。

输入文件

1. 药物相关死亡（DRDs）数据集：

- 文件名： `refined_data/allDRDs.rds`
- 用途：这是主要的输入数据集，包含了药物相关死亡的详细信息。

2. 全因死亡（ACM）数据集：

- 文件名： `lookup_NRS_all_cause_mortality.rds` 和 `lookup_NRS_all_cause_mortality_ADDITIONAL_CASES.rds`
- 用途：这两个文件包含了全因死亡的数据。代码中将这两个文件合并后，用于进一步的分析和比较。

输出文件

1. 处理后的DRDs数据集：

- 文件名： `refined_data/allDRDs.rds`
- 内容：经过变量重命名、筛选和转换后的DRDs数据集。这个文件包含了财年、年龄等新计算的变量。

2. 死亡计数数据集：

- 文件名： `refined_data/death_counts.csv` （注释掉了，未实际生成）

- 内容：包含不同财年、不同条件下的死亡计数。例如，特定财年内的总死亡数、特定类型的死亡数等。

3. 处理后的ACM数据集：

- 文件名： `refined_data/ACM_drds.rds`
- 内容：经过变量重命名、筛选和转换后的ACM数据集。这个文件包含了财年、年龄等新计算的变量。

4. 所有HMMB死亡数据集：

- 文件名： `refined_data/allHMMBs.rds`
- 内容：筛选出的与HMMB（可能是指某种特定药物或药物类别）相关的死亡数据。这个数据集经过了详细的筛选和标记，包括特定的死亡类型、年龄范围、药物标记等。

5. 所有ACM数据集：

- 文件名： `refined_data/allACMs.rds`
- 内容：经过处理后的全因死亡数据集，包含了标记和筛选后的数据。

代码

Code block

```
1  # Create cohort; alive & with prescriptions in 2011-2020
2  # Last updated 30/08/2022
3
4  # Libraries and data -----
5  rm(list=ls())
6  library(tidyverse)
7  library(tidylog)
8  library(janitor)
9  library(lubridate)
10 library(retroharmonize)
11 library(data.table)
12 library(phsmethods)
13
14
15 #load DRDs file
16 setwd(' ')
17 getwd()
18 drds <- read_rds('refined_data/allDRDs.rds')
```

```

19 #load secondary codes file
20
21 # 指定路径读取药物相关死亡 (DRDs) 数据集。重命名变量
22 drds <- drds %>% rename(SMR99_HMMBimplic = NRSdrd_HMMBimplic) %>%
23 rename(SMR99_DOB = NRSdrd_DOB) %>%
24 rename(SMR99_Sex = NRSdrd_SEX) %>%
25 rename(SMR99_DOD = NRSdrd_DOD) %>%
26 rename(SMR99_HBres = NRSdrd_HBres) %>%
27 rename(SMR99_UNDERLYING_CAUSE_OF_DEATH = NRSdrd_UNDERLYING_CAUSE_OF_DEATH) %>%
28 rename(type_of_death = type_of_death) %>%
29 rename(SMR99_poison = NRSdrd_poison) %>%
30 rename(SMR99_alsopres = NRSdrd_alsopres)
31
32 # 转换变量 选择特定的列。
33 # 将IAIN列转换为字符类型。
34 # 使用phsmethods::extract_fin_year函数提取财年。
35 # 计算年龄。
36 drds <- drds %>%
37   select(IAIN, SMR99_HMMBimplic, SMR99_DOB, SMR99_Sex, SMR99_DOD,
38         SMR99_HBres, SMR99_UNDERLYING_CAUSE_OF_DEATH, type_of_death,
39         SMR99_poison, SMR99_alsopres, strong_opioid_flag) %>%
40   mutate(across(c("IAIN"), as.character)) %>%
41   mutate(FY := phsmethods::extract_fin_year(lubridate::ymd(SMR99_DOD)))
42   %>%
43   mutate(age= floor(as.duration(SMR99_DOB %--% SMR99_DOD) /
44     ddays(365.25)))
45
46 # 检查数据的范围和唯一值
47 range(drds$SMR99_DOD)
48 length(unique(drds$IAIN))
49 #7385
50 length(unique(drds$IAIN[which(drds$SMR99_HMMBimplic==1)]))
51 #5514 (HMMBs)
52 length(unique(drds$IAIN[which(drds$SMR99_HMMBimplic==1 &
53   drds$type_of_death==3)]))
54 #4831 (accidental HMMBs)
55
56 #write_rds(drds, "refined_data/allDRDs.rds")
57 #####
58 #####
59
60 # 计算特定财年内的死亡计数。
61 death_counts <- drds %>%
62   filter(SMR99_HMMBimplic==1) %>%

```

```

63     filter(FY %in% c('2014/15', '2015/16', '2016/17', '2017/18',
64                     '2018/19', '2019/20', '2020/21', '2021/22',
65                     '2022/23')) %>%
66     count(FY)
67
68
69 death_counts <- death_counts %>%
70   rename(n0=n)
71
72
73 death_counts2 <- drds %>%
74   filter(SMR99_HMMBimplic==1) %>%
75   filter(FY %in% c('2014/15', '2015/16', '2016/17', '2017/18',
76                   '2018/19', '2019/20', '2020/21', '2021/22',
77                   '2022/23')) %>%
78   filter(type_of_death==3) %>%
79   count(FY)
80
81 # 将不同条件下的死亡计数合并到一个数据框中
82 death_counts <- death_counts %>%
83   left_join(death_counts2 %>% rename(n2=n))
84
85 death_counts2 <- drds %>%
86   filter(SMR99_HMMBimplic==1) %>%
87   filter(FY %in% c('2014/15', '2015/16', '2016/17', '2017/18',
88                   '2018/19', '2019/20', '2020/21', '2021/22',
89                   '2022/23')) %>%
90   filter(type_of_death==3) %>%
91   filter(age>=15 & age<65) %>%
92   # filter(age>65 & age<70) %>%
93   count(FY)
94
95 death_counts <- death_counts %>%
96   left_join(death_counts2 %>% rename(n3=n))
97
98
99
100
101 death_counts2 <- drds %>%
102   filter(SMR99_HMMBimplic==1) %>%
103   filter(FY %in% c('2014/15', '2015/16', '2016/17', '2017/18',
104                   '2018/19', '2019/20', '2020/21', '2021/22',
105                   '2022/23')) %>%
106   filter(type_of_death==3) %>%
107   filter(age>=15 & age<65) %>%
108   # filter(age>65 & age<70) %>%
109   filter(SMR99_UNDERLYING_CAUSE_OF_DEATH %in% c('X42', 'F112', 'F192')) %>%

```



```

110     count(FY)
111
112 death_counts <- death_counts %>%
113     left_join(death_counts2 %>% rename(n4=n))
114
115
116
117 death_counts2 <- drds %>%
118     filter(SMR99_HMMBimplic==1) %>%
119     filter(FY %in% c('2014/15', '2015/16', '2016/17', '2017/18',
120                     '2018/19', '2019/20', '2020/21', '2021/22',
121                     '2022/23')) %>%
122     filter(type_of_death==3) %>%
123     filter(age>=15 & age<65) %>%
124     # filter(age>65 & age<70) %>%
125     filter(SMR99_UNDERLYING_CAUSE_OF_DEATH %in% c('X42', 'F112', 'F192')) %>%
126     filter(strong_opioid_flag==0) %>%
127     count(FY)
128
129
130 death_counts <- death_counts %>%
131     left_join(death_counts2 %>% rename(n5=n))
132
133
134 #write.csv(death_counts, "refined_data/death_counts.csv")
135
136 #####
137 #####
138 #####
139 #####
140 #####
141 #####
142 #####
143 #####
144 #####
145 #####
146 #####
147 #####
148 # add ACM ----
149
150
151 ACM <- read_rds("lookup_NRS_all_cause_mortality.rds")
152 ExtraACM <- read_rds("lookup_NRS_all_cause_mortality_ADDITIONAL_CASES.rds")
153 colnames(ACM)==colnames(ExtraACM)
154 ACM <- rbind(ACM,ExtraACM)
155
156 colnames(ACM) <- c('rowid', 'IAIN', 'SMR99_DOB', 'SMR99_DOD', 'SMR99_Sex',

```

```

157         'record_type', 'date_of_registration',
158         'SMR99_UNDERLYING_CAUSE_OF_DEATH',
159         'SMR99_CAUSE_OF_DEATH_CODE_0',
160         'SMR99_CAUSE_OF_DEATH_CODE_1',
161         'SMR99_CAUSE_OF_DEATH_CODE_2',
162         'SMR99_CAUSE_OF_DEATH_CODE_3',
163         'SMR99_CAUSE_OF_DEATH_CODE_4',
164         'SMR99_CAUSE_OF_DEATH_CODE_5',
165         'SMR99_CAUSE_OF_DEATH_CODE_6',
166         'SMR99_CAUSE_OF_DEATH_CODE_7',
167         'SMR99_CAUSE_OF_DEATH_CODE_8',
168         'SMR99_CAUSE_OF_DEATH_CODE_9',
169         'SMR99_HBres', 'SMR99_Council_Area',
170         'SMR99_HBres_current', 'SMR99_Council_Area_current',
171         'SMR992_flag')
172
173
174
175
176
177 ACM <- ACM %>%
178   select(IAIN, SMR99_DOB, SMR99_Sex, SMR99_DOD,
179          SMR99_HBres, SMR99_UNDERLYING_CAUSE_OF_DEATH, ) %>%
180   mutate(across(c("IAIN"), as.character)) %>%
181   mutate(FY := phsmethods::extract_fin_year(lubridate::ymd(SMR99_DOD))) %>%
182   mutate(age_at_death = floor(as.duration(SMR99_DOB %--% SMR99_DOD) /
183   ddays(365.25)))
184
185
186 # Check overlap between ACM and DRDs 重叠?
187
188 length(unique(ACM$IAIN))
189
190 length(intersect(unique(ACM$IAIN), unique(drds$IAIN)))
191
192
193 #write_rds(ACM, "refined_data/ACM_drds.rds")
194 #筛选年龄在15-64岁、有具体原因且无强效阿片类药物的HMMB意外死亡
195
196 all_HMMB_deaths <- drds[which(drds$SMR99_HMMBimplic==1),]
197
198
199
200 allHMMBs <- all_HMMB_deaths %>%
201   filter(type_of_death==3) %>%
202   mutate(FY := phsmethods::extract_fin_year(lubridate::ymd(SMR99_DOD))) %>%

```

```

203   mutate(age_at_death = floor(as.duration(SMR99_DOB %--% SMR99_DOD) /
days(365.25))) %>%
204   filter(age_at_death>=15 & age_at_death<=64) %>%
205   mutate(type_of_death_long = case_when(
206     type_of_death==3 ~ "Accidental",
207     type_of_death==1 ~ "Suicide",
208     type_of_death==2 ~ "Undetermined")) %>%
209   mutate(age_group = case_when(age_at_death>=15 & age_at_death<35 ~ '15-34',
210                                age_at_death>=35 & age_at_death<50 ~ '35-49',
211                                age_at_death>=50 & age_at_death<65 ~ '50-64'))
%>%
212   mutate(across(c("IAIN"), as.character)) %>%
213   mutate(sex= ifelse(SMR99_Sex==1, 'M', 'F')) %>%
214   filter(strong_opioid_flag==0) %>%
215   filter(SMR99_UNDERLYING_CAUSE_OF_DEATH %in% c('X42', 'F112', 'F192')) %>%
216   mutate(day=SMR99_DOD) %>%
217   mutate(DOB=SMR99_DOB) %>%
218   mutate(hmmb_flag=1) %>%
219   filter(type_of_death==3)
220
221
222
223   #allHMMBs <- allHMMBs %>%
224   #filter(FY %in% c('2014/15', '2015/16', '2016/17',
225   #                  '2017/18', '2018/19', '2019/20',
226   #                  '2020/21', '2021/22', '2022/23'))
227   #allHMMBs <- allHMMBs %>% filter(strong_opioid_flag==0)
228
229
230   # Prepare all-cause mortality dataset for export
231
232   all_deaths <- read_rds("refined_data/ACM_drds.rds")
233   all_deaths <- all_deaths %>%
234     mutate(across(c("IAIN"), as.character))
235
236   allACM <- all_deaths %>%
237     mutate(DOB=SMR99_DOB) %>%
238     mutate(day=SMR99_DOD) %>%
239     mutate(acm_flag=1) %>%
240     select(IAIN, day, DOB, acm_flag)
241
242
243   write_rds(allHMMBs, "refined_data/allHMMBs.rds")
244   write_rds(allACM, "refined_data/allACMs.rds")
245
246
247

```

00d-Hospital_Admissions

输入文件

- `SMR0104_drug_related_hosp_stays.rds`：这是一个RDS文件，包含与药物相关的医院住院数据。该文件被加载到 `hospital_admissions` 变量中。

输出文件

- `refined_data/hospitalisations_T.rds`：这是一个RDS文件，包含经过筛选和处理后的与特定ICD-10代码（T400, T401, T403）相关的住院数据。该文件存储在 `refined_data` 目录中。

数据加载和预处理

1. 加载必要的库：

2. 设置工作目录：

- 使用 `setwd('')` 设置工作目录（需要用户指定路径）。
- 使用 `getwd()` 确认当前工作目录。

3. 加载数据：

- 使用 `read_rds` 函数加载 `SMR0104_drug_related_hosp_stays.rds` 文件到 `hospital_admissions` 变量中。

4. 数据预处理：

- 检查 `ADMISSION_DATE` 和 `flag_discharge_death` 字段的范围。
- 将 `IAIN` 字段转换为字符类型。
- 提取财政年度（`FY`）。
- 计算入院时的年龄（`age at entry`）。

- 过滤掉 `flag_discharge_death` 为1的记录（即死亡记录）。
- 过滤年龄在15到64岁之间的记录。
- 去重（`distinct`）。

数据筛选和统计

1. 统计财政年度的记录数：

- 使用 `count(FY)` 统计每个财政年度的记录数。

2. 筛选特定ICD-10代码的记录：

- 筛选 `RECORD_TYPE` 为 `01B` 的记录。
- 筛选主诊断或次要诊断为 `T400` , `T401` , `T403` 的记录。
- 过滤年龄在15到64岁之间的记录。
- 过滤掉 `flag_discharge_death` 为1的记录。

3. 统计性别分布：

- 对2014/15到2021/22财政年度的性别分布进行统计。
- 对2014/15到2022/23财政年度的性别分布进行统计。

数据保存

1. 保存处理后的数据：

- 使用 `write_rds` 函数将处理后的数据保存到 `refined_data/hospitalisations_T.rds` 文件中。

总结

这段代码的主要功能是：

1. 加载和预处理与药物相关的医院住院数据。
2. 筛选出特定ICD-10代码（T400, T401, T403）的记录。
3. 统计不同财政年度和性别的记录数。
4. 保存处理后的数据到指定文件中。

代码

Code block

```
1  # Split dataset by age_grp and year_grp
2  # Calculate CMRs for on/off by year_grp, and by age_grp and year_grp
3  # Plot figures
4
5  # Libraries and data -----
6  rm(list=ls())
7  library(tidyverse)
8  library(tidylog)
9  library(lubridate)
10 library(glue)
11 library(data.table)
12 library(janitor)
13 #library(epitools)
14 #library(haven)
15 library(xtable)
16 library(phsmethods)
17
18
19 setwd('')
20 getwd()
21 hospital_admissions <- read_rds('SMR0104_drug_related_hosp_stays.rds')
22
23
24 # Check date ranges and discharge death flag
25
26 range(hospital_admissions$ADMISSION_DATE)
27 range(hospital_admissions$flag_discharge_death)
28
29
30 # Prepare dataset: convert ID to character, extract financial year, calculate
  age at entry
31
32 hospital_admissions <- hospital_admissions %>%
33   mutate(across(c("IAIN"), as.character)) %>%
34   mutate(FY := phsmethods::extract_fin_year(lubridate::ymd(ADMISSION_DATE)))
  %>%
35   mutate(age_at_entry = floor(as.duration(DOB %--% ADMISSION_DATE) /
  ddays(365.25))) %>%
36   filter(flag_discharge_death==0)
37   #mutate(age_at_entry = round(as.duration(DOB %--% ADMISSION_DATE) /
  ddays(365.25)))
```

```

38
39
40 # Count admissions by financial year
41
42 hospital_admissions %>% count(FY)
43
44
45 # Filter for working-age adults (15-64) and remove duplicates
46
47 hospital_admissions <- hospital_admissions %>%
48     filter(age_at_entry>=15 & age_at_entry<=64) %>%
49     filter(flag_discharge_death==0)
50
51
52 # Count by record type with totals
53
54 hospital_admissions %>% count(RECORD_TYPE) %>% adorn_totals()
55 hospital_admissions <- hospital_admissions %>% distinct()
56
57 #length(which(hospital_admissions$EPI_1_MAIN_CONDITION=='F112'))
58
59
60 # Count T-code admissions by sex for selected financial years
61
62 hosp_T_codes <- hospital_admissions %>%
63     filter(RECORD_TYPE %in% c('01B')) %>%
64     filter(EPI_1_MAIN_CONDITION %in% c('T400','T401','T403') |
65           EPI_1_OTHER_CONDITION_1 %in% c('T400','T401','T403') |
66           EPI_1_OTHER_CONDITION_2 %in% c('T400','T401','T403') |
67           EPI_1_OTHER_CONDITION_3 %in% c('T400','T401','T403') |
68           EPI_1_OTHER_CONDITION_4 %in% c('T400','T401','T403') |
69           EPI_1_OTHER_CONDITION_5 %in% c('T400','T401','T403')) %>%
70     filter(age_at_entry>=15 & age_at_entry<=64) %>%
71     filter(flag_discharge_death==0)
72
73 hosp_T_codes %>%
74     filter(FY %in% c('2014/15','2015/16','2016/17','2017/18',
75                     '2018/19','2019/20','2020/21','2021/22')) %>%
76     count(SEX) %>% adorn_totals()
77
78 hosp_T_codes %>%
79     filter(FY %in% c('2014/15','2015/16','2016/17','2017/18',
80                     '2018/19','2019/20','2020/21','2021/22','2022/23')) %>%
81     count(SEX) %>% adorn_totals()
82
83 write_rds(hosp_T_codes, "refined_data/hospitalisations_T.rds")
84

```

00e_create_demographics

输入文件

- `/OST_prescriptions.rds`：这是一个RDS文件，包含阿片类药物替代治疗（ORT）的处方数据。该文件被加载到 `orts` 变量中。

输出文件

- `refined_data/demographics.rds`：这是一个RDS文件，包含经过处理和合并后的人口统计学数据（包括性别、出生日期和NHS健康委员会信息）。该文件存储在 `refined_data` 目录中。

代码

2 读入原始处方数据（ORT）

Code block

```
1  setwd('') # 设置项目工作目录（已留空，需手动指定）
2  orts <- readRDS(.../filtered_ORT_prescriptions.rds) # 读入 2011-2020 的 ORT 处方
3  orts <- orts %>% mutate(across("IAIN", as.character)) # 保证 ID 为字符型
```

- **注意：** `setwd('')` 需改成实际路径，否则后续读写会报错。

3 把日期列统一转成 `Date`

Code block

```
1  z_orts <- orts %>%
2    mutate_at(vars(contains("Date")), ~as.Date(., "%Y/%m/%d %H:%M:%S"))
3  z_orts <- z_orts %>% mutate(IAIN = as.character(IAIN))
```


- **目的：**所有含 `Date` 的列（如 `PaidDate`、`PrescriptionStartDate` 等）从字符转成真正的日期，便于后续按时间筛选与排序。

4 Health Board (HB) 名称转换函数

Code block

```
1 hb_recode <- function(hb_in, recode_to = "long", and_choice = "and") { ... }
```

- **功能：**把 NHS Health Board 的代码（如 `S08000007`）或简写（如 `GC`）映射为 ****完整名称****（`Greater Glasgow and Clyde`）或 ****简写****（`GC`）。
- **参数：**
 - `recode_to = "long"` → 返回完整名称（如 `"Greater Glasgow and Clyde"`）
 - `recode_to = "short"` → 返回两字母简写（如 `"GC"`）
- **注意：**函数里用 `case_when` 做了大量手工映射，保证新旧编码兼容。

5 提取人口学变量 (Gender + DOB)

Code block

```
1 demographics <- orts %>% select(IAIN, Gender, DOB)
2 demographics <- demographics %>% tidylog::distinct() # 去重
3 demographics <- demographics %>%
4   mutate(DOB = as.Date(DOB, "%Y/%m/%d %H:%M:%S")) %>%
5   mutate(IAIN = as.character(IAIN))
```

- **目的：**把 `Gender`（1=M, 2=F）和出生日期单独存一份，后面与死亡、HB 合并。

6 为每个患者每年生成“最近居住 HB”

6-a) 先拿 按年+人 的最近记录

Code block

```
1 NHSboards2 <- orts %>%
2   mutate(year = year(PaidDate)) %>%
3   select(IAIN, PatientsResidentsNHSBoard, PaidDate, year) %>%
4   filter(!is.na(PatientsResidentsNHSBoard)) %>%
5   group_by(IAIN, year) %>%
6   arrange(desc(PaidDate)) %>% # 同一年里按处方日期倒序
7   slice(1) %>% # 取最近一条处方
8   ungroup() %>% distinct()
```

- **结果：**得到 **每人每年** 的最新 HB 代码（用于后续时间序列插补）。

6-b) 再拿 全时期 的最近 HB（备用）

Code block

```
1 NHSboards <- orts %>%
2   select(IAIN, PatientsResidentsNHSBoard, PaidDate) %>%
3   filter(!is.na(...)) %>%
4   group_by(IAIN) %>% arrange(desc(PaidDate)) %>% slice(1) %>%
5   ungroup() %>% mutate(year = year(PaidDate)) %>% distinct()
```

- **用途：**做稳健性检查，也可作为“缺失年份”的 fallback。

7 把 HB 缺失年份 向前/向后填充

Code block

```
1 n_years <- length(2009:2024) # 16 年
2 n_IAIN <- length(unique(NHSboards2$IAIN))
3
4 NHSBoards_by_year <- expand_grid(
5   IAIN = unique(NHSboards2$IAIN),
6   year = 2009:2024
```

```

7 ) %>% left_join(NHSboards2) # 先左连，很多缺失
8
9 # 用 locf 两次：先向后、再向前填充
10 NHSBoards_by_year2 <- NHSBoards_by_year %>%
11   arrange(IAIN, year) %>%
12   group_by(IAIN) %>%
13   mutate(PatientsResidentsNHSBoard = zoo::na.locf(PatientsResidentsNHSBoard,
14     na.rm = FALSE, fromLast = TRUE)) %>%
14   mutate(PatientsResidentsNHSBoard = zoo::na.locf(PatientsResidentsNHSBoard,
15     na.rm = FALSE, fromLast = FALSE)) %>%
15   ungroup()

```

- 目的：即使某年没有处方，也能用 **最近已知 HB** 填充，确保 2009–2024 每年每人都有值。

8 合并人口学 + HB，并保存

Code block

```

1 demographics_and_NHSBoard2 <- NHSBoards_by_year2 %>%
2   select(IAIN, PatientsResidentsNHSBoard, year) %>%
3   left_join(demographics) %>% distinct()
4
5 demographics_and_NHSBoard <- demographics %>%
6   left_join(NHSboards %>% select(IAIN, PatientsResidentsNHSBoard))
7
8 # 写盘
9 write_rds(demographics_and_NHSBoard, "refined_data/demographics.rds")

```

- 最终产物： `demographics.rds` 包含 每人每年的
 - IAIN
 - Gender
 - DOB
 - 最近居住 NHS Board（已插补）

