

Format WA Data - 2016 Estimates

Jeanette Birnbaum

9/26/2017

Contents

1	Raw Data Overview	1
1.1	Sample Size	1
1.2	Variable list	2
1.3	Variable summaries	2
2	Subset based on hst=WA and year	6
2.1	First, split the combined year-quarter of diagnosis and AIDS variables	6
2.2	Subset the data based on hst=WA and year	6
2.3	New sample size	7
3	Year and quarter of diagnosis: cleaning it up	7
3.1	Years represented	7
3.2	Quarters represented	7
3.3	Distribute unknown quarters uniformly across Q1-Q4	8
4	Tabulate and collapse race and mode of diagnosis variables	8
4.1	Race and mode by year	8
4.2	Collapse	9
5	AIDS at Diagnosis	10
5.1	AIDS at initial diagnosis?	10
5.2	Years of AIDS diagnosis represented:	10
5.3	Quarters of AIDS diagnosis represented:	10
6	Ever had a last negative test (everHadNegTest)	10
6.1	Coding	10
6.2	Make compatible with recorded LNT dates	11
7	Time since last negative test (infPeriod)	12
7.1	Apply age-16 assumption and summarize	12
7.2	Diagnoses younger than 16	12
7.3	Maximum window of 18 years	13
8	Final analytic dataset	14
8.1	Reminder of data cleaning	14
8.2	Variable summaries	14

1 Raw Data Overview

1.1 Sample Size

N = 21098

1.2 Variable list

```
str(dataf)

## 'data.frame':    21098 obs. of  21 variables:
## $ firstvl      : num  8433 19914 35382 51 108 ...
## $ firstcd4cnt  : num  177 243 501 636 847 ...
## $ tth_ever_neg : int   5 5 5 5 5 5 5 5 5 5 ...
## $ new_race     : Factor w/ 8 levels "White","Black",...: 2 2 1 1 1 1 3 1 1 1 ...
## $ hst          : chr   "WA" "WA" "WA" "WA" ...
## $ hdx_age      : int   51 25 41 34 38 33 33 41 45 19 ...
## $ new_mode     : Factor w/ 9 levels "MSM","IDU","MSM/IDU",...: 3 6 6 1 1 1 3 1 1 1 ...
## $ tth_lneg_dt_flag : int  4 4 4 4 4 4 4 4 4 4 ...
## $ tth_ppos_dt_flag : int  4 4 4 4 4 4 4 4 4 4 ...
## $ est_infect_period: int  3 3 3 3 3 3 3 3 3 3 ...
## $ hdx_yr_qtr    : chr   "1998_3Q" "1999_3Q" "1995_2Q" "1990_" ...
## $ hdx_dt_flag   : chr   "M" "M" "M" "Y" ...
## $ adx_yr_qtr    : chr   "2003_2Q" "2000_1Q" NA NA ...
## $ adx_dt_flag   : chr   "M" "M" NA NA ...
## $ lag_lneg_hdx_dt : int  NA NA NA NA NA NA NA NA NA NA ...
## $ lag_ppos_hdx_dt : int  NA NA NA NA NA NA NA NA NA NA ...
## $ tth_prev_pos  : chr   "N" "N" "N" "N" ...
## $ dx_in_king    : chr   "Y" "Y" "Y" "Y" ...
## $ vl_days       : int  673 111 7517 4032 1396 3061 2618 1810 1607 4461 ...
## $ cd4_days      : int  1734 122 7517 6294 4151 3857 2618 2283 2350 5356 ...
## $ meth_use      : chr   NA NA NA NA ...
```

1.3 Variable summaries

```
##
##
##
## VARIABLE 1 : firstvl
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
##      0.0   921.5 17400.0 38720.0 98870.0 100000.0   6756
##
##      Percent missing:[1] 32.02
##
##
##
## VARIABLE 2 : firstcd4cnt
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
##      0.0   150.0   344.0   387.9   564.0  4269.0   6581
##
##      Percent missing:[1] 31.19
##
##
##
## VARIABLE 3 : tth_ever_neg
##      var
##      1    2    5 <NA>
## 3031  736 17331    0
##
```

```

##      Percent missing:[1] 0
##
##
##
## VARIABLE 4 : new_race
##      var
##      White   Black   Hisp   Asian   NHoPI   AI/AN   Multi   Unknown   <NA>
##      14368   2826   2167   565     78       299     785     10         0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 5 : hst
##      var
##      WA  <NA>
##      21098      0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 6 : hdx_age
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00   28.00   35.00   35.82   42.00   91.00
##
##      Percent missing:[1] NA
##
##
##
## VARIABLE 7 : new_mode
##      var
##      MSM      IDU      MSM/IDU      Transfus      Hemo
##      13308     1689     1986         122         101
##      Hetero     Ped F Pres Hetero      NIR      <NA>
##      2025       118         0         1749         0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 8 : tth_lneg_dt_flag
##      var
##      1      2      3      4  <NA>
##      471  1700   768 18159      0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 9 : tth_ppos_dt_flag
##      var
##      1      2      3      4  <NA>
##      1121  2332   345 17300      0

```

```

##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 10 : est_infect_period
##      var
##      1      2      3 <NA>
## 1605 1022 18471      0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 11 : hdx_yr_qtr
##      [1] ""
##
##      Percent missing:numeric(0)
##
##
##
## VARIABLE 12 : hdx_dt_flag
##      var
##      D      M      Y <NA>
## 5712 13105 2281      0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 13 : adx_yr_qtr
##      [1] ""
##
##      Percent missing:numeric(0)
##
##
##
## VARIABLE 14 : adx_dt_flag
##      var
##      D      M      Y <NA>
## 2951 11670      58 6419
##
##      Percent missing:[1] 30.42
##
##
##
## VARIABLE 15 : lag_lneg_hdx_dt
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
##      0.0   176.5   438.0   958.3  1138.0 11570.0  18159
##
##      Percent missing:[1] 86.07
##
##
##

```

```

## VARIABLE 16 : lag_ppos_hdx_dt
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##      0.0      0.0      5.0   342.6    13.0 13380.0   17300
##
##      Percent missing:[1] 82
##
##
##
## VARIABLE 17 : tth_prev_pos
##      var
##      N      Y  <NA>
## 20417   681     0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 18 : dx_in_king
##      var
##      N      Y  <NA>
##  8052 13046     0
##
##      Percent missing:[1] 0
##
##
##
## VARIABLE 19 : vl_days
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##      0         7      64    1288    2301   11440   6701
##
##      Percent missing:[1] 31.76
##
##
##
## VARIABLE 20 : cd4_days
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##      0         7     114    1481    2679   11440   6573
##
##      Percent missing:[1] 31.15
##
##
##
## VARIABLE 21 : meth_use
##      var
##      NO UNKNOWN    YES    <NA>
##      703      203    417   19775
##
##      Percent missing:[1] 93.73

```

2 Subset based on hst=WA and year

2.1 First, split the combined year-quarter of diagnosis and AIDS variables

```
#####  
# SPLIT COMBINED YR-QTR VARIABLE  
#####  
# Year, quarter, and quarter-year of Dx (diagnosis)  
dataf$yearDx <- as.numeric(substring(dataf$hdx_yr_qtr,0,4))  
dataf$quarterDx <- as.numeric(substring(dataf$hdx_yr_qtr,6,6))  
dataf$timeDx <- dataf$yearDx + (dataf$quarterDx-1)/4  
# AIDS at Dx - if missing, assumed to be false  
dataf$aidsAtDx <- dataf$hdx_yr_qtr == dataf$adx_yr_qtr  
dataf$aidsAtDx[is.na(dataf$aidsAtDx)] <- FALSE  
# Year, quarter, and quarter-year of AIDS (diagnosis)  
dataf$yearAids <- as.numeric(substring(dataf$adx_yr_qtr,0,4))  
dataf$quarterAids <- as.numeric(substring(dataf$adx_yr_qtr,6,6))  
dataf$timeAids <- dataf$yearAids + (dataf$quarterAids-1)/4
```

2.2 Subset the data based on hst=WA and year

```
#####  
# SUBSET THE DATA - INITIAL RESTRICTIONS  
#####  
if (!'year_min'%in%ls()) year_min <- 2005  
if (!'year_max'%in%ls()) year_max <- 2013  
  
# Year min and max for this run  
c(year_min, year_max)  
  
## [1] 2005 2016  
  
# Non-sequential look  
table(hst_included=dataf$hst=='WA', useNA='ifany')  
  
## hst_included  
## TRUE  
## 21098  
  
table(yearDx_included=dataf$yearDx>=year_min & dataf$yearDx<=year_max,  
      useNA='ifany')  
  
## yearDx_included  
## FALSE TRUE  
## 14987 6111  
  
table(yearDx_missing=is.na(dataf$hdx_yr_qtr))  
  
## yearDx_missing  
## FALSE  
## 21098  
  
table(age_missing_and_missing_lastNeg=(is.na(dataf$hdx_age) &  
                                         is.na(dataf$lag_lneg_hdx_dt)))
```

```

## age_missing_and_missing_lastNeg
## FALSE
## 21098

# Sequential look
(hst_included <- table(hst_included=dataf$hst=='WA', useNA='ifany'))

## hst_included
## TRUE
## 21098

dataf <- subset(dataf, hst=='WA')
(yearDx_included <- table(yearDx_included=(dataf$yearDx>=year_min & dataf$yearDx<=year_max), useNA='ifany'))

## yearDx_included
## FALSE TRUE
## 14987 6111

dataf <- subset(dataf, yearDx>=year_min & yearDx<=year_max)
(age_included <- table(age_and_lastNeg_present=! (is.na(dataf$hdx_age) &
                                                    is.na(dataf$lag_lneg_hdx_dt))))

## age_and_lastNeg_present
## TRUE
## 6111

dataf <- subset(dataf, !(is.na(hdx_age) & is.na(lag_lneg_hdx_dt)))
(Nobs1 <- nrow(dataf))

## [1] 6111

```

Excluded 14987 cases based on year and hst restrictions and missingness in age and year of diagnosis.

2.3 New sample size

New sample size is 6111

3 Year and quarter of diagnosis: cleaning it up

3.1 Years represented

```

##
## 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
## 561 535 583 539 546 556 493 509 455 444 458 432

```

3.2 Quarters represented

```

##
## 1 2 3 4 <NA>
## 1589 1562 1513 1436 11

```

3.3 Distribute unknown quarters uniformly across Q1-Q4

```
#####
# IMPUTE A QUARTER IF ONLY YEAR IS KNOWN
#####
impute_qtr <- !is.na(dataf$yearDx) & is.na(dataf$quarterDx)
set.seed(98103)
dataf$quarterDx[impute_qtr] <- sample(4, size=sum(impute_qtr),
                                     replace=TRUE)
dataf$timeDx <- dataf$yearDx + (dataf$quarterDx-1)/4
summary(dataf$timeDx, digits=6)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 2005.00 2007.50 2010.25 2010.56 2013.50 2016.75

time_min <- min(dataf$timeDx)
time_max <- max(dataf$timeDx)

# Time min and max for this run
c(time_min, time_max)

## [1] 2005.00 2016.75
```

4 Tabulate and collapse race and mode of diagnosis variables

4.1 Race and mode by year

```
table(dataf$new_race, dataf$yearDx, useNA='ifany')
```

```
##
##           2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
## White      315  314  326  283  317  319  280  284  243  227  224  202
## Black       99   77  100  102   91   80   89   95   88   96   92   91
## Hisp        76   65   90   95   86  104   77   64   78   61   86   71
## Asian       20   23   22   28   25   26   24   30   24   38   35   36
## NHOPI        2    5    2    0    3    1    5    7    5    5    4    4
## AI/AN        7    7    6   11    6    8    4    5    4    6    5   10
## Multi       42   44   37   20   18   18   14   24   13   11   12   18
## Unknown     0    0    0    0    0    0    0    0    0    0    0    0

table(dataf$new_mode, dataf$yearDx, useNA='ifany')
```

```
##
##           2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015
## MSM          295  311  335  301  316  351  297  282  266  251  274
## IDU           40   42   32   25   27   33   30   22   20   23   36
## MSM/IDU       65   43   50   34   43   30   47   42   33   29   22
## Transfus       1    0    1    1    0    0    0    0    0    0    0
## Hemo           1    0    0    0    0    0    0    0    0    0    0
## Hetero        92   70   82   85   74   68   39   40   36   35   36
## Ped           0    3    2    3   11   10    6    3    4    3    4
## F Pres Hetero  0    0    0    0    0    0    0    0    0    0    0
## NIR           67   66   81   90   75   64   74  120   96  103   86
```



```
##
##           2016
##   MSM           224
##   IDU           29
##   MSM/IDU       24
##   Transfus      0
##   Hemo          0
##   Hetero        47
##   Ped           5
##   F Pres Hetero 0
##   NIR           103
```

4.2 Collapse

```
#####
# COLLAPSE RACE AND MODE OF DIAGNOSIS
#####

race_levels <- c('White', 'Black', 'Hisp', 'Asian', 'Native', 'Multi')
mode_levels <- c('MSM', 'Hetero', 'Blood/Needle')
dataf <- within(dataf, {
  race <- as.character(new_race)
  race[race=='AI/AN' | race == 'NHoPI'] <- 'Native'
  race <- factor(race,
                 labels=race_levels,
                 levels=race_levels)
  mode <- as.character(new_mode)
  mode[mode=='MSM/IDU'] <- 'MSM'
  mode[mode=='F Pres Hetero' | mode=='NIR'] <- 'Hetero'
  mode[mode=='IDU' | mode=='Transfus' | mode=='Hemo' |
        mode=='Ped'] <- 'Blood/Needle'
  mode <- factor(mode,
                 levels=mode_levels,
                 labels=mode_levels)
  mode2 <- factor(ifelse(mode=='MSM', 'MSM', 'non-MSM'))
})
```

```
table(dataf$race, dataf$yearDx, useNA='ifany')
```

```
##
##           2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
##   White    315  314  326  283  317  319  280  284  243  227  224  202
##   Black     99   77  100  102   91   80   89   95   88   96   92   91
##   Hisp      76   65   90   95   86  104   77   64   78   61   86   71
##   Asian     20   23   22   28   25   26   24   30   24   38   35   36
##   Native     9   12    8   11    9    9    9   12    9   11    9   14
##   Multi     42   44   37   20   18   18   14   24   13   11   12   18
```

```
table(dataf$mode, dataf$yearDx, useNA='ifany')
```

```
##
##           2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
##   MSM           360  354  385  335  359  381  344  324  299  280  296  248
##   Hetero        159  136  163  175  149  132  113  160  132  138  122  150
```

```
## Blood/Needle 42 45 35 29 38 43 36 25 24 26 40 34
```

```
table(dataf$mode2, dataf$yearDx, useNA='ifany')
```

```
##
##           2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
## MSM       360 354 385 335 359 381 344 324 299 280 296 248
## non-MSM   201 181 198 204 187 175 149 185 156 164 162 184
```

5 AIDS at Diagnosis

5.1 AIDS at initial diagnosis?

```
##
## FALSE TRUE
## 4555 1556
```

5.2 Years of AIDS diagnosis represented:

```
##
## 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 <NA>
## 166 208 219 265 273 237 233 200 175 159 166 137 34 3639
```

5.3 Quarters of AIDS diagnosis represented:

```
##
## 1 2 3 4 <NA>
## 611 648 616 592 3644
```

6 Ever had a last negative test (everHadNegTest)

6.1 Coding

This variable will be coded as Yes=TRUE, No=FALSE, and Don't Know/Refused/Missing=NA

```
#####
# CREATE everHadNegTest
#####
# Define everHadNegTest based on tth_ever_neg
# 2015 data update: this variable was coded numerically, so I have
# added that option in.
dataf <- transform(dataf,
                    everHadNegTest=ifelse(tth_ever_neg=='Y' | tth_ever_neg==1, TRUE,
                                           ifelse(tth_ever_neg=='N' | tth_ever_neg==2, FALSE, NA)))
with(dataf, table(everHadNegTest, tth_ever_neg, useNA='always'))
```

```
##           tth_ever_neg
## everHadNegTest 1 2 5 <NA>
## FALSE 0 711 0 0
## TRUE 2941 0 0 0
```

```
##           <NA>      0      0 2459      0
# Now cross-check it with the lag_lneg_hdx_dt, which actually has the
# time since last negative test
(checkEver <- with(dataf, table(everHadNegTest,
                                TID_NA=is.na(lag_lneg_hdx_dt), useNA='always'))))

##           TID_NA
## everHadNegTest FALSE TRUE <NA>
##           FALSE      7  704      0
##           TRUE    2830  111      0
##           <NA>     18 2441      0

# Look at actual lag_lneg_hdx_dt values by everHadNegTest
ddply(dataf, .(everHadNegTest), function(x) c(summary(x$lag_lneg_hdx_dt)))

##   everHadNegTest Min. 1st Qu. Median  Mean 3rd Qu. Max. NA's
## 1             FALSE  101   354.0   681.0 644.9    975 1074   704
## 2              TRUE    0   176.0   434.0 952.5   1134 9925   111
## 3              NA   122   208.8   467.5 803.5   1399 2476  2441
```

6.2 Make compatible with recorded LNT dates

6.2.1 Change incorrect FALSEs

We have 7 cases with everHadNegTest=FALSE and 18 with everHadNegTest=NA but have a time since last negative test. Change their everHadNegTest flag.

```
toTRUE1 <- !dataf$everHadNegTest & !is.na(dataf$lag_lneg_hdx_dt)
toTRUE2 <- is.na(dataf$everHadNegTest) & !is.na(dataf$lag_lneg_hdx_dt)
dataf$everHadNegTest[toTRUE1] <- TRUE
dataf$everHadNegTest[toTRUE2] <- TRUE
```

6.2.2 Change incorrect TRUEs

We have 111 cases who have everHadNegTest=TRUE but have NO time since last negative test. Change their everHadNegTest flag. Change, 9/27/17 - previously was setting to false; now, set to NA.

```
## an alternative to setting to FALSE
toNA <- dataf$everHadNegTest & is.na(dataf$lag_lneg_hdx_dt)
dataf$everHadNegTest[toNA] <- NA
```

6.2.3 Check

```
(checkEver <- with(dataf, table(everHadNegTest,
                                TID_NA=is.na(lag_lneg_hdx_dt), useNA='always'))))

##           TID_NA
## everHadNegTest FALSE TRUE <NA>
##           FALSE      0  704      0
##           TRUE    2855    0      0
##           <NA>      0 2552      0
```

7 Time since last negative test (infPeriod)

7.1 Apply age-16 assumption and summarize

```
#####
# CREATE infPeriod and then look at it
#####

#### TEMPORARY:
#dataf$age=35

aidsUB <- qweibull(.95,shape=2.516,scale=1/0.086) #17.98418
dataf <- within(dataf,{
  lastNeg_yrs=lag_lneg_hdx_dt/365
  infPeriod=ifelse(everHadNegTest,
    pmin(lastNeg_yrs, aidsUB),
    ifelse(!everHadNegTest,
      pmin(hdx_age-16, aidsUB),
      NA))
  earliestInf=hdx_age-infPeriod
})

summary(dataf$infPeriod,digits=3)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
## -3.000   0.601   1.940   4.940   6.940   18.000    2552
```

7.2 Diagnoses younger than 16

```
# Number of cases who got a negative infPeriod
(neginfPeriod <- sum(dataf$infPeriod<0,na.rm=TRUE))

## [1] 4

# Diagnoses at or under age 16 by everHadNegTest
(a1 <- table(atunder16=dataf$hdx_age<=16,
  everHadNegTest=dataf$everHadNegTest, useNA='ifany'))

##           everHadNegTest
## atunder16 FALSE TRUE <NA>
##      FALSE   696 2850 2474
##      TRUE     8    5   78

# Diagnoses at or under age 16 by year, 2005-2013
table(atunder16count=subset(dataf, yearDx>=year_min & yearDx<=year_max)$hdx_age<=16,
  year=subset(dataf, yearDx>=year_min & yearDx<=year_max)$yearDx, useNA='ifany')

##           year
## atunder16count 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
##      FALSE   558  529  577  532  534  544  485  497  445  438  453  428
##      TRUE     3    6    6    7   12   12    8   12   10    6    5    4

# Now just under 16, excluding hdx_age=16
# Diagnoses under age 16 by everHadNegTest
```

```
(a2 <- table(under16=dataf$hdx_age<16,
             everHadNegTest=dataf$everHadNegTest, useNA='ifany'))

##           everHadNegTest
## under16 FALSE TRUE <NA>
##   FALSE   700 2852 2478
##   TRUE     4    3   74

# Diagnoses under age 16 by year
table(under16count=subset(dataf, yearDx>=year_min & yearDx>=year_max)$hdx_age<16,
      year=subset(dataf, yearDx>=year_min & yearDx>=year_max)$yearDx, useNA='ifany')

##           year
## under16count 2016
##   FALSE   428
##   TRUE     4

# Among those diagnosed at or under 16: everHadNegTest by mode
table(everHadNegTest=subset(dataf,hdx_age<=16)$everHadNegTest,
      mode=subset(dataf,hdx_age<=16)$new_mode, useNA='ifany')
```

```
##           mode
## everHadNegTest MSM IDU MSM/IDU Transfus Hemo Hetero Ped F Pres Hetero NIR
##   FALSE     3   1      0      0   0      1   3      0   0
##   TRUE      1   1      1      0   0      1   0      0   1
##   <NA>      2   0      0      0   0      1  45      0  30
```

There are 86 cases who do not have a date of last negative test and may not fit the assumption of TID=age-16. Of those, 8 are age 16 at diagnosis and will have TID=0 using this assumption. Primary mode of transmission is Ped ('Perinatal or pediatric').

```
(young_included <- with(dataf,
                        table(over16_or_atunder16_with_obs_infPeriod=
                              (hdx_age>16 |
                               !(hdx_age<=16 & (!everHadNegTest |
                                                    is.na(everHadNegTest)))))))

## over16_or_atunder16_with_obs_infPeriod
## FALSE TRUE
##   86 6025

dataf <- subset(dataf, !(hdx_age<=16 & (!everHadNegTest |
                                         is.na(everHadNegTest))))

(Nobs2 <- nrow(dataf))

## [1] 6025

summary(dataf$infPeriod, digits=3)
```

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
## 0.000  0.607   1.950   4.950   6.980  18.000   2474
```

Excluded 86 cases due to age ≤ 16 and no observed infPeriod data.

7.3 Maximum window of 18 years

```
# We did cap some people whose TID's were >aidsUB
(check_cap1 <- with(subset(dataf, everHadNegTest),
  table(original_over_aidsUB=lastNeg_yrs>aidsUB,
    infPeriod_over_aidsUB=infPeriod>aidsUB,
    useNA='ifany')))
```

```
##                infPeriod_over_aidsUB
## original_over_aidsUB FALSE
##                FALSE  2824
##                TRUE    31
```

Among those with everHadNegTest=TRUE, we capped 31 cases at aidsUB.

```
(check_cap2 <- with(subset(dataf, !everHadNegTest),
  table(original_over_aidsUB=lastNeg_yrs>aidsUB,
    infPeriod_over_aidsUB=infPeriod>aidsUB,
    useNA='ifany')))
```

```
##                infPeriod_over_aidsUB
## original_over_aidsUB FALSE
##                <NA>    696
```

Among those with everHadNegTest=FALSE, no one had an original TID value.

```
(check_cap3 <- with(subset(dataf, is.na(everHadNegTest)),
  table(original_over_aidsUB=lastNeg_yrs>aidsUB,
    infPeriod_over_aidsUB=infPeriod>aidsUB,
    useNA='ifany')))
```

```
##                infPeriod_over_aidsUB
## original_over_aidsUB <NA>
##                <NA>  2474
```

Among those with everHadNegTest=NA, no one had an original TID value.

8 Final analytic dataset

8.1 Reminder of data cleaning

Final subset is of size 6025 * Diagnoses included: - Year: non-missing, and 2005 onwards - Occurred in WA state - Excluded 14987 cases based on year and hst restrictions (no missingness in age and year of diagnosis in data for 2015 estimates). * Ages included: - If missing age, must have recorded time of last negative test - If age ≤ 16 , must have recorded time of last negative test - Excluded 86 cases due to age ≤ 16 and no observed LNT.

8.2 Variable summaries

```
## [1] 6025
##
## VARIABLE: hdx_age
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   14.00  28.00   36.00   37.46   46.00   83.00
##
## VARIABLE: timeDx
```

```

##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##      2005      2008      2010      2011      2014      2017
##
## VARIABLE: everHadNegTest
##      Mode  FALSE    TRUE    NA's
## logical    696    2855    2474
##
## VARIABLE: lastNeg_yrs
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.    NA's
##      0.000  0.482   1.189   2.605   3.096   27.190   3170
##
## VARIABLE: infPeriod
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.    NA's
##      0.0000  0.6068  1.9450  4.9510  6.9750  17.9800   2474

```