

PIRLS-test.R

Carl Schmertmann

Thu Aug 22 16:09:22 2019

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 3.5.3
```

```
## -- Attaching packages -----
```

```
## v ggplot2 3.1.0      v purrr   0.3.0
## v tibble  2.0.1      v dplyr  0.8.0.1
## v tidyr   0.8.2      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.3.0
```

```
## Warning: package 'tibble' was built under R version 3.5.2
```

```
## Warning: package 'readr' was built under R version 3.5.2
```

```
## Warning: package 'purrr' was built under R version 3.5.2
```

```
## Warning: package 'dplyr' was built under R version 3.5.2
```

```
## Warning: package 'stringr' was built under R version 3.5.2
```

```
## -- Conflicts -----
```

```
## x dplyr::filter() masks stats::filter()
```

```
## x dplyr::lag()     masks stats::lag()
```

```
library(broom)
```

```
## Warning: package 'broom' was built under R version 3.5.3
```

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

```
##### DATA #####
```

```
## Italy 1980 Female data from HMD (true e0 from HMD is 77.42)
```

```
ITA = read.csv(file='ITA-Female-1980.csv')
```

```
# standard schedule = smoothed CAN females 1959 log rates at 0,1,...99
```

```
# read the std schedule (log rates for CAN females 1959, ages 0...99)
```

```
std = read.csv('female-std.csv')$std
```

```
#####
```

```
# note that this sources TOPALS_fit.R (the grouped version)
```

```
# rather than TOPALS_fit function.R (the single-year version)
```

```
source('TOPALS_fit.R')
```

```
# some utility functions
```

```
## age-grouping function
```

```
agg = function(x,bounds) {
```

```
  age = seq(x)-1 # 0,1,2,...
```

```

L = head(bounds,-1)
U = tail(bounds,-1)
as.vector( tapply( x, cut(age, breaks=bounds, right=FALSE), sum))
}

## plotting function
show_fit = function(fit, true_schedule, fit_color='red') {

  df_grouped = data.frame(
    L = head(fit$age_group_bounds,-1),
    U = tail(fit$age_group_bounds,-1),
    N = fit$N,
    D = fit$D
  ) %>%
    mutate(logmx_obs = log(D/N))

  df_single = data.frame(
    age      = seq(std)-0.5,
    std      = myfit$std,
    logmx_true = true_schedule,
    logmx_fit  = myfit$logm
  )

  this_plot =
    ggplot(data = df_single, aes(x=age,y=logmx_true)) +
    geom_line(aes(x=age,y=std), color='black', lwd=0.2) +
    geom_line(aes(x=age,y=logmx_fit), color=fit_color, lwd=3, alpha=.40) +
    geom_segment(data=df_grouped,
      aes(x=L,xend=U,y=logmx_obs,yend=logmx_obs),
      color=fit_color,lwd=1, alpha=.90) +
    geom_point(size=0.80) +
    labs(x='Age',y='Log Mortality Rate',
      title='Italy Females 1980',
      subtitle = paste(sum(D),'deaths to',round(sum(N)),'women')) +
    scale_x_continuous(breaks=c(0,1,seq(5,100,5)),minor_breaks = NULL) +
    theme_bw()

  print(this_plot)
} # show_fit

# trapez approx of life expectancy from a logmx schedule over ages 0..99
e0 = function(logmx) {
  mx = exp(logmx)
  px = exp(-mx)
  lx = c(1,cumprod(px))
  return( sum(head(lx,-1) + tail(lx,-1)) / 2)
}

#####
# FULL DATASET WITH 1-YEAR GROUPS
#####

```

```
## full dataset: 0,1,...99
N = ITA$N[1:100]
D = ITA$D[1:100]

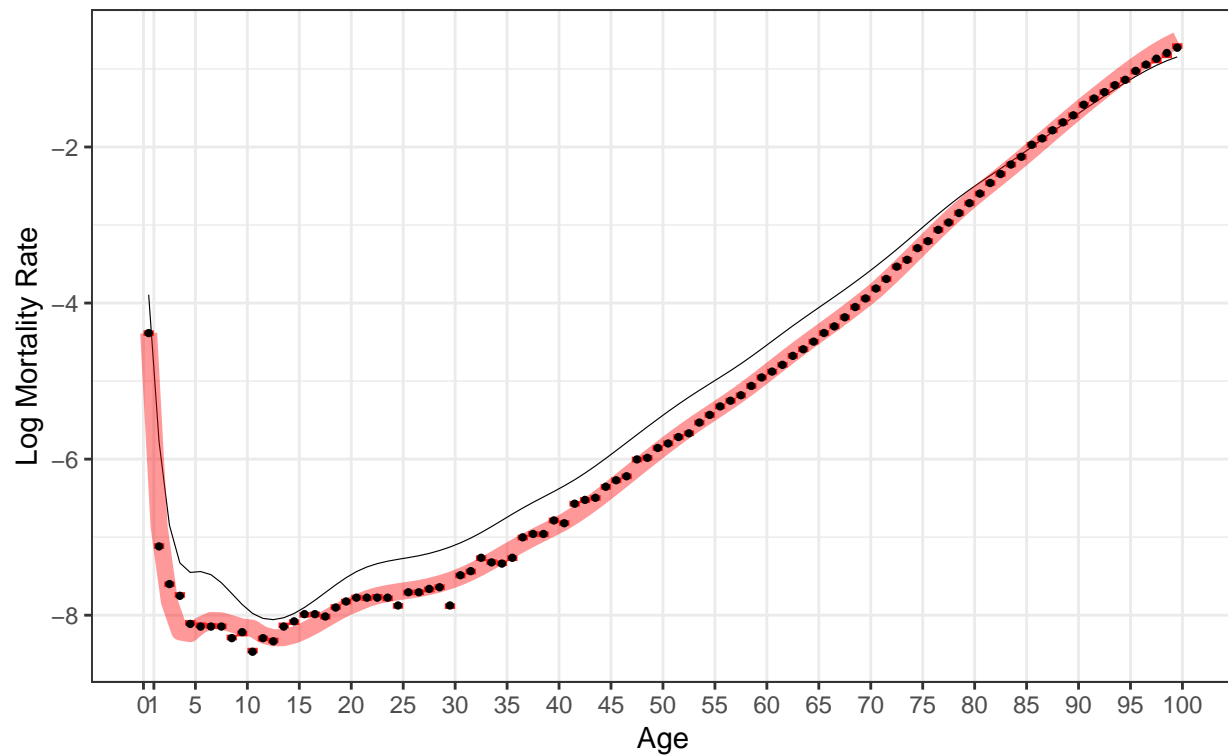
bounds = 0:100

myfit = TOPALS_fit(D=D, N=N, std=std,
                  age_group_bounds = bounds,
                  details=TRUE)

## Loading required package: splines
show_fit(myfit, true_schedule = ITA$logmx[1:100],
         fit_color = 'red')
```

Italy Females 1980

259667 deaths to 28950657 women



```
#####
# FULL DATASET WITH 5-YEAR GROUPS
#####

bounds = c(0,1,seq(5,100,5))

N = agg(ITA$N, bounds)
D = agg(ITA$D, bounds)

myfit = TOPALS_fit(D=D, N=N, std=std,
                  age_group_bounds = bounds,
```

```

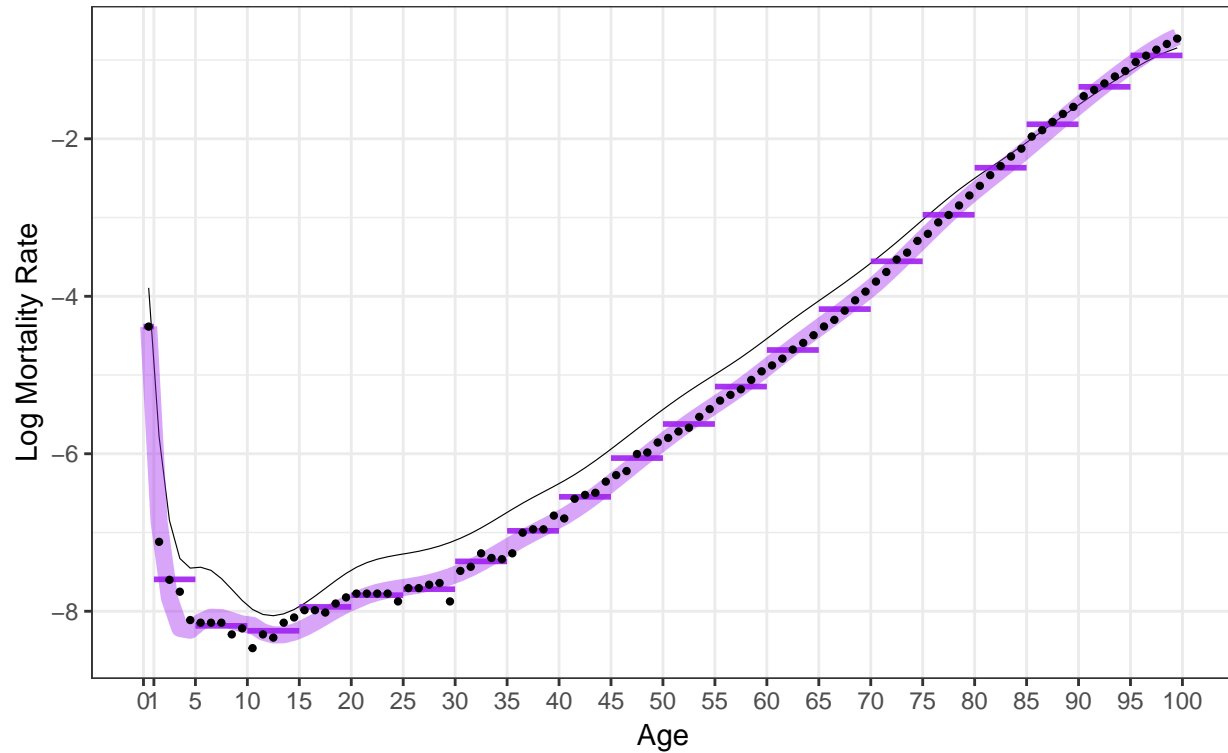
                                details=TRUE)

show_fit( myfit,  true_schedule = ITA$logmx[1:100],
          fit_color = 'purple')

```

Italy Females 1980

259667 deaths to 28950657 women



```

#####
# FULL DATASET WITH 10-YEAR GROUPS
#####

bounds = c(0,1, seq(10,100,10))

N = agg(ITA$N, bounds)
D = agg(ITA$D, bounds)

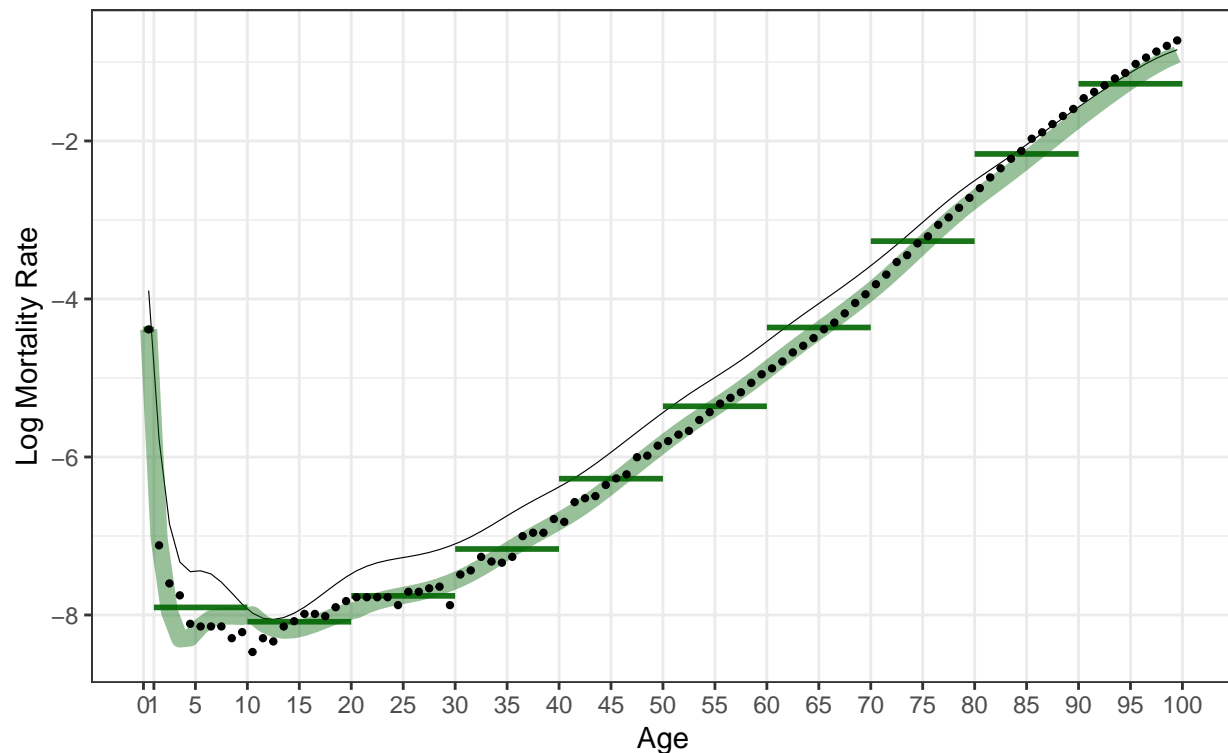
myfit = TOPALS_fit(D=D, N=N, std=std,
                  age_group_bounds = bounds,
                  details=TRUE)

show_fit( myfit,  true_schedule = ITA$logmx[1:100],
          fit_color = 'darkgreen')

```

Italy Females 1980

259667 deaths to 28950657 women



SMALL POPULATION SIMULATIONS

scenario is a data frame with experimental parameters
target_pop, L, U, nsim on each row

```
nsim = 500
```

```
pop_vals = c(5e5, 1e5, 1e4, 1e3)
```

```
bnd_vals = list( seq(0,100,1),
                  c(0,1,seq(5,100,5)),
                  c(0,1,seq(10,100,10))
                )
```

MAE and density variables for each scenario will be calculated below

```
scenario = expand.grid( target_pop = pop_vals,
                      bounds      = bnd_vals,
                      MAE         = Inf) %>%
  as_tibble()
```

add an empty LIST column to hold densities

```
scenario = scenario %>%
  add_column(e0_dens = list(NA))
```

```
for (s in 1:nrow(scenario)) {
```

```

target_pop = scenario$target_pop[s]
bounds      = unlist( scenario$bounds[s] )

## exposure and deaths for these age groups (all of Italy)
bigN = agg(ITA$N, bounds)
bigD = agg(ITA$D, bounds)

e      = rep(NA,nsim)

# small population with same age structure as ITA
N      = bigN * target_pop/sum(bigN)

for (i in 1:nsim) {

  # random deaths for this small popualtion at Italian rates
  D = rpois(length(N), N * bigD/bigN)

  myfit = TOPALS_fit(D=D, N=N, std=std,
                    age_group_bounds = bounds,
                    details=TRUE)

  e[i] = e0(myfit$logm)
} # for i

scenario$MAE[s]      = round(mean( abs(e-77.42)),2)
scenario$e0_dens[s] = list( tidy(density(e, adj=1.5) ))

} # for s

# MAE report
matrix( round(scenario$MAE,2), nrow=4,
        dimnames=list(paste('Pop=',format(pop_vals,scientific = FALSE)),
                      paste0(c('1','5','10'),'-yr grp')))

##           1-yr grp 5-yr grp 10-yr grp
## Pop= 500000    0.17    0.20    0.52
## Pop= 100000    0.37    0.37    0.56
## Pop=  10000    1.06    1.06    1.19
## Pop=   1000    3.40    3.15    3.14

## e0 densities

for (p in unique(scenario$target_pop)) {

  tmp = filter( scenario, target_pop==p)

  df1 = as.data.frame(tmp$e0_dens[1]) %>%
    add_column(grouping=1)
  df2 = as.data.frame(tmp$e0_dens[2]) %>%
    add_column(grouping=5)
  df3 = as.data.frame(tmp$e0_dens[3]) %>%
    add_column(grouping=10)

```

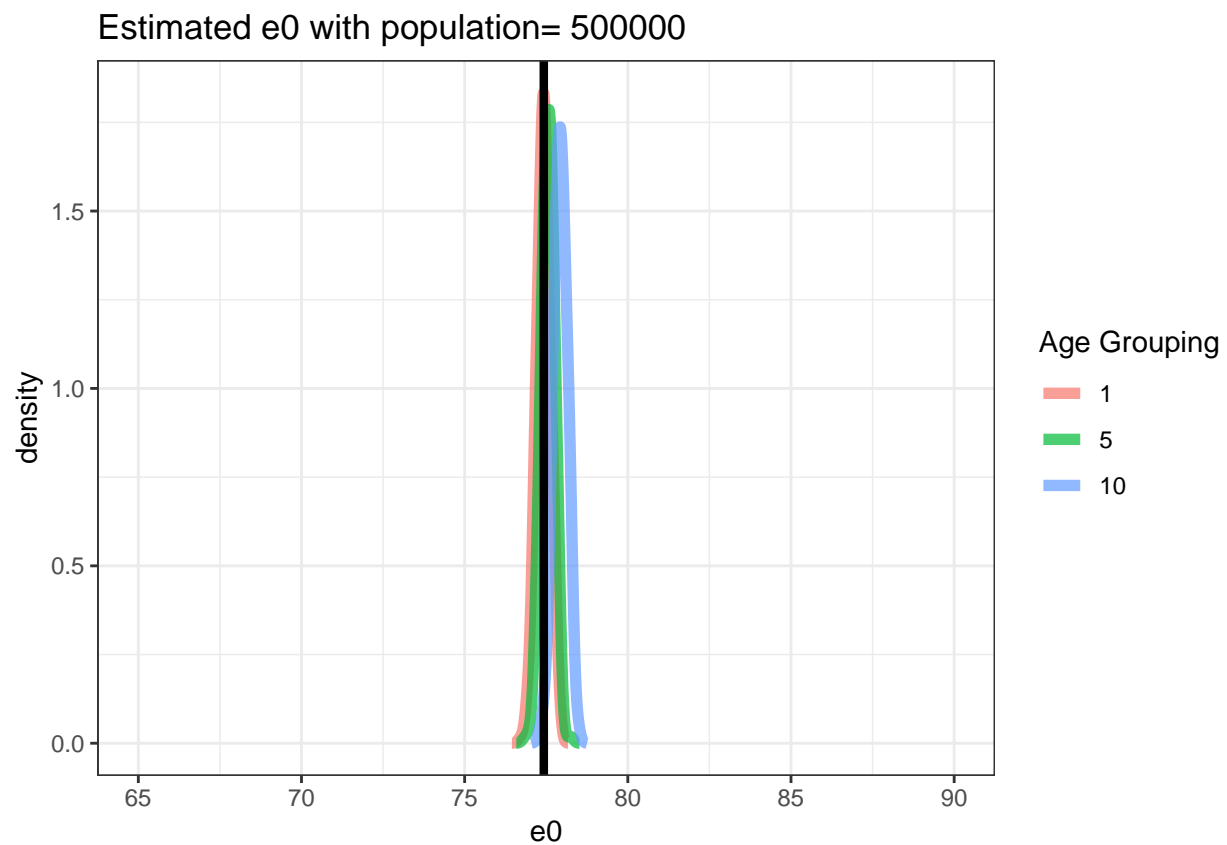
```

df = rbind(df1,df2,df3)

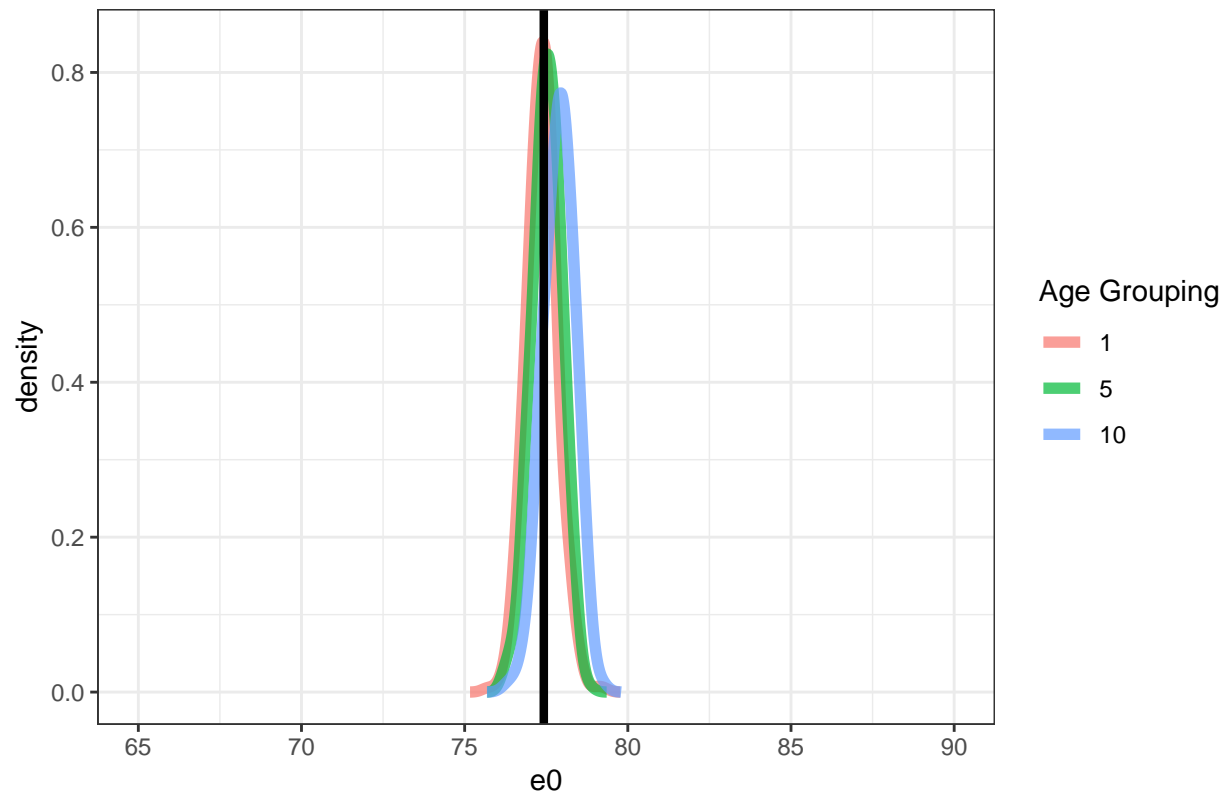
this_plot =
  ggplot(data=df, aes(x=x,y=y, color=as.factor(grouping))) +
    geom_line(lwd=2, alpha=.70) +
    labs(title=paste('Estimated e0 with population=',
                     format(p,scientific = FALSE)),
         x='e0',y='density', color='Age Grouping') +
    geom_vline(xintercept = e0(ITA$logmx), lwd=1.5) +
    scale_x_continuous(limits=c(65,90)) +
    theme_bw()

print(this_plot)
}

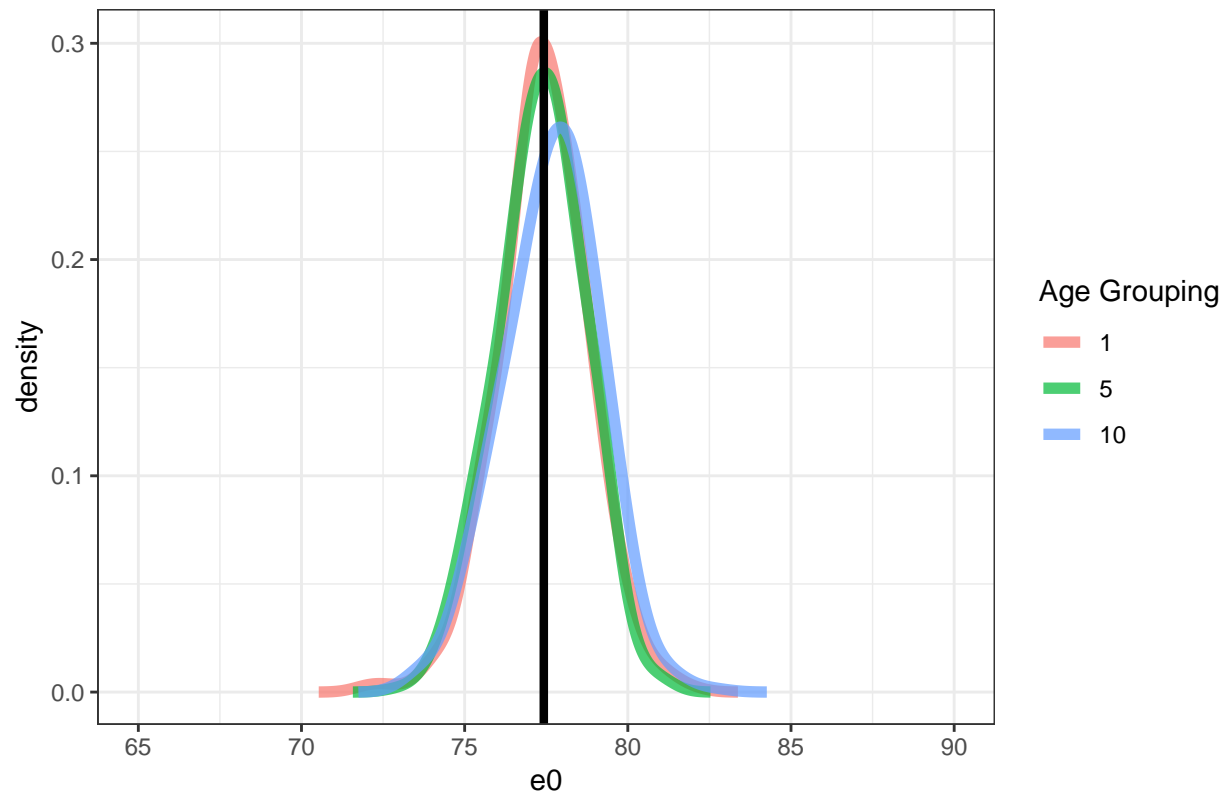
```



Estimated e0 with population= 100000



Estimated e0 with population= 10000



Warning: Removed 576 rows containing missing values (geom_path).

