

# Bias correction in GLMs

Bicko, Jonathan & Ben

2021 Jun 21 (Mon)

## Introduction

We intend to investigate our prediction based on known truth and any bias potentially introduced by non-linear averaging, conditioning or random effect. We'll start with a simple case of a only fixed effect model and then consider a mixed effect model.

## Simulation

We perform a simple simulation for a fixed effect model

$$\begin{aligned}\text{logit}(\text{status} = 1) &= \eta \\ \eta &= \beta_0 + \beta_A \text{Age} + \beta_W \text{Wealthindex} \\ \text{Age} &\sim \text{Normal}(0.2, 1) \\ \text{Wealthindex} &\sim \text{Normal}(0, 1) \\ \beta_0 &= 0.7 \\ \beta_A &= 0.3 \\ \beta_W &= 0.6\end{aligned}$$

```
N <- 1e5
beta0 <- 0.7
betaA <- 0.2
betaW <- 0.5

age_max <- 1
age_min <- 0.2
# age <- runif(N, age_min, age_max)
age <- rnorm(N, age_min, age_max)

wealthindex <- rnorm(N, 0, 1)

eta <- beta0 + betaA * age + betaW * wealthindex
sim_df <- (data.frame(age=age, wealthindex=wealthindex, eta=eta)
  %>% mutate(status = rbinom(N, 1, plogis(eta)))
  %>% select(-eta)
)
true_prop <- mean(sim_df$status)
print(true_prop)
```

```
## [1] 0.70092
```

```
head(sim_df)
```

```
##           age wealthindex status
## 1 1.8654889   0.9422537      1
## 2 0.3995828  -0.5698525      0
## 3 0.2830708  -0.2246213      0
## 4 1.6965558  -2.8186033      1
## 5 2.2311610  -1.6049063      0
## 6 1.2664415   0.5309227      1
```

## Simple logistic model

```
simple_mod <- glm(status ~ age + wealthindex, data = sim_df, family="binomial")
```

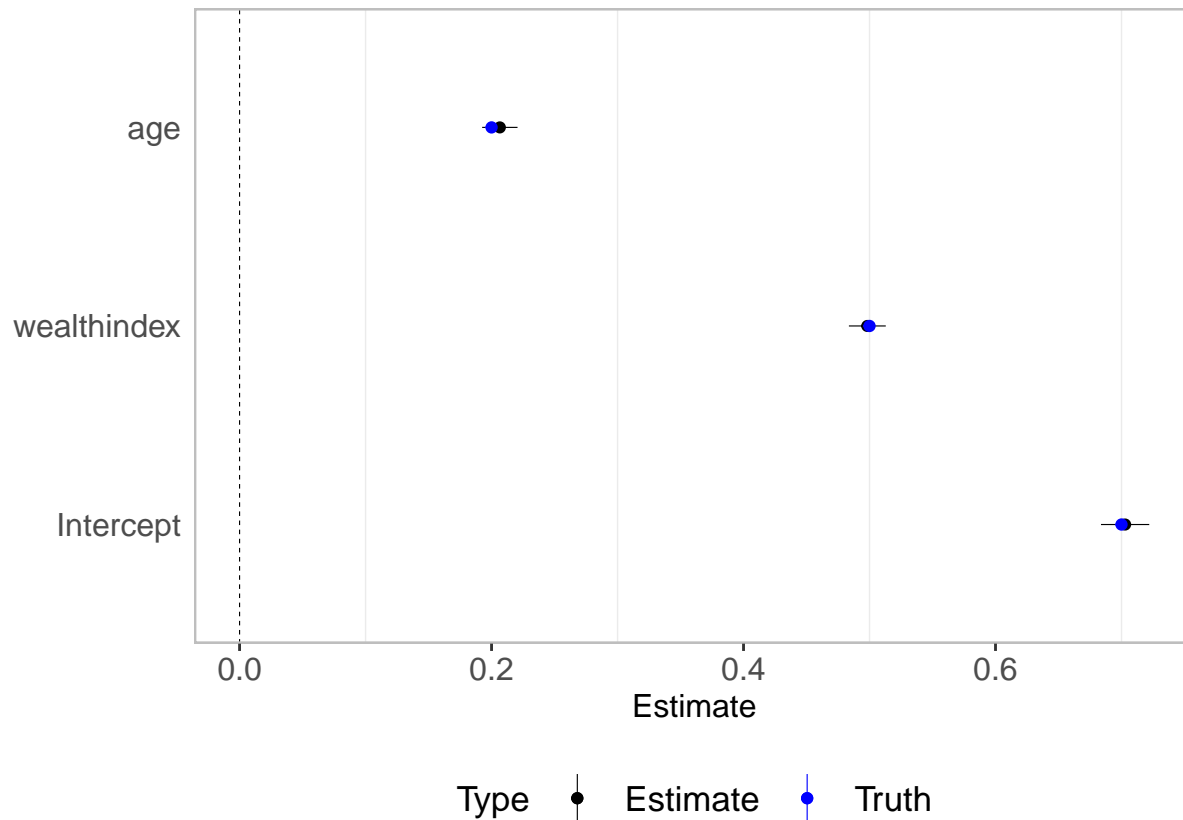
Coefficient plots

```
## True beta
true_beta_df <- data.frame(term=c("Intercept", "age", "wealthindex")
, estimate=c(beta0, betaA, betaW)
)

## Tidy coef estimates
coef_df <- (broom::tidy(simple_mod, conf.int=TRUE)
# %>% dotwhisker::by_2sd(sim_df)
# %>% mutate(term = gsub("\\(|\\)", "", term))
)
print(coef_df)

## # A tibble: 3 x 7
##   term          estimate std.error statistic    p.value conf.low conf.high
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>    <dbl>    <dbl>
## 1 Intercept      0.703     0.00977     71.9 0.    0.684     0.722
## 2 age            0.207     0.00716     28.9 4.55e-183 0.193     0.221
## 3 wealthindex    0.498     0.00747     66.7 0.    0.484     0.513

simple_coef_plot <- (plotEsize(coef_df)
+ geom_point(data=true_beta_df, aes(x=term, y=estimate, colour="Truth"))
+ labs(colour="Type")
+ scale_colour_manual(values=c("black", "blue"))
)
print(simple_coef_plot)
```



Variable effect plots – varpred and population averaging approach

```
# Age
## varpred way
simple_vareff_age <- varpred(simple_mod, "age", isolate=FALSE, modelname="varpred")

# Wealth index
## varpred - no bias adjustment
simple_vareff_wealthindex <- varpred(simple_mod, "wealthindex", isolate=TRUE, modelname="varpred")

## Bias adjusted
simple_vareff_wealthindex_adjust <- varpred(simple_mod, "wealthindex", isolate=TRUE, pop.ave=TRUE, modelname="varpred")

vareff_wealthindex <- simple_vareff_wealthindex
vareff_wealthindex$preds <- do.call("rbind", list(vareff_wealthindex$preds, simple_vareff_wealthindex_adjust$preds))
wealthindex_plot <- (plot(vareff_wealthindex)
  + labs(y="", colour="Model")
  + geom_hline(yintercept=true_prop, lty=2, colour="grey")
  + scale_colour_manual(values=c("black", "blue"))
  + theme(legend.position="bottom")
)

## Scale for 'colour' is already present. Adding another scale for 'colour',
## which will replace the existing scale.
```

Effect sizes on logit scale:

```
coef_df_logit <- (coef_df
  %>% select(term, estimate, conf.low, conf.high)
```

```

    %>% group_by(term)
    %>% summarise_all(plogis)
  )
print(coef_df_logit)

```

```

## # A tibble: 3 x 4
##   term          estimate conf.low conf.high
##   <chr>          <dbl>    <dbl>    <dbl>
## 1 age            0.551      0.548      0.555
## 2 Intercept      0.669      0.665      0.673
## 3 wealthindex    0.622      0.619      0.625

```

## Variable predictions

We consider both **varpred** and population averaging approach; and then introduce bias correction to **varpred** predictions.

```

popavefun <- function(mod, focal, non.focal, level=0.95, modelname="Pop. ave", ...) {
  mf <- model.matrix(mod)
  mm <- (mf
    %>% data.frame()
    %>% mutate_at(non.focal, mean)
    %>% as.matrix()
  )
  vc <- vcov(mod)
  linpred <- as.vector(mm %*% coef(mod))
  pse_var <- sqrt(rowSums(mm * t(tcrossprod(data.matrix(vc), mm))))
  z.val <- qnorm(1 - (1 - level)/2)
  pred_df <- (mf
    %>% data.frame()
    %>% select_at(focal)
    %>% mutate(fit = linpred
      , lwr = plogis(fit - z.val*pse_var)
      , upr = plogis(fit + z.val*pse_var)
      , fit = plogis(fit)
      , model = modelname
      , se = NA
    )
  )
  return(pred_df)
}

```

- Age

```

## varpred way
simple_vareff_age <- varpred(simple_mod, "age", isolate=FALSE, modelname="varpred")

## Pop. average
simple_vareff_age_pop <- popavefun(simple_mod, "age", "wealthindex", modelname = "Pop. ave")

## Bias adjust
simple_vareff_age_adjust <- varpred(simple_mod, "age", isolate=TRUE, bias.adjust=TRUE, modelname = "Bias adj")

vareff_age <- simple_vareff_age
vareff_age$preds <- do.call("rbind", list(vareff_age$preds, simple_vareff_age_pop, simple_vareff_age_ad

```

```
age_plot <- (plot(vareff_age)
+ labs(y="Prob. of improved \n service", colour="Model")
+ geom_hline(yintercept=true_prop, lty=2, colour="grey")
+ scale_colour_manual(values=c("black", "blue", "red"))
+ theme(legend.position="bottom")
)
```

## Scale for 'colour' is already present. Adding another scale for 'colour',  
## which will replace the existing scale.

- Wealth index

```
# Wealth index
## varpred
simple_vareff_wealthindex <- varpred(simple_mod, "wealthindex", isolate=TRUE, modelname="varpred")

## Pop. average
simple_vareff_wealthindex_pop <- varpred(simple_mod, "wealthindex", isolate=TRUE, pop.ave=TRUE, modelname="pop.ave")

## Bias adjust
simple_vareff_wealthindex_adjust <- varpred(simple_mod, "wealthindex", isolate=TRUE, bias.adjust=TRUE, modelname="bias.adjust")

vareff_wealthindex <- simple_vareff_wealthindex
vareff_wealthindex$preds <- do.call("rbind", list(vareff_wealthindex$preds, simple_vareff_wealthindex_pop$preds, simple_vareff_wealthindex_adjust$preds))
wealthindex_plot <- (plot(vareff_wealthindex)
+ labs(y="", colour="Model")
+ geom_hline(yintercept=true_prop, lty=2, colour="grey")
+ scale_colour_manual(values=c("black", "blue", "red"))
+ theme(legend.position="bottom")
)
```

## Scale for 'colour' is already present. Adding another scale for 'colour',  
## which will replace the existing scale.

```
ggarrange(age_plot, wealthindex_plot, common.legend=TRUE)
```

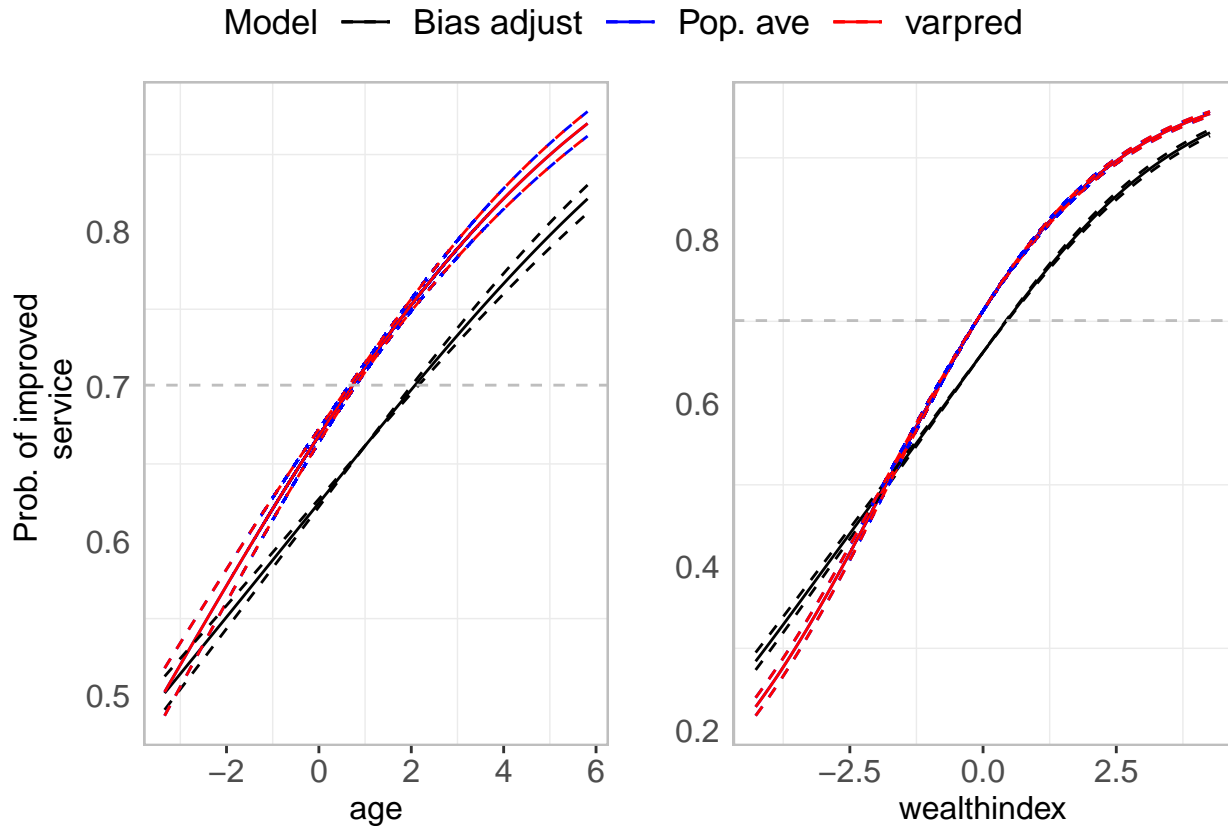


Figure 1: A comparison of population averaged and varpred-based predictions. We also apply bias-adjustment to the ‘vapred’ predictions. For ‘age’, we implement the naive approach to compute the predictions in ‘popavefun’ function and then implement the same in ‘vareffects’ so as to use the centering machineries. In both cases, the population averaging and varpred gives similar estimates. The estimated population average is very close to the observed in the case of ‘wealthindex’ but slightly higher in the case of ‘age’ for the unadjusted estimates (but very low for bias-adjusted). , see previous paragraph.

The observed population average is 0.70092 while the estimated population averages (similar estimates with varpred) are:

- age: unadjusted 0.7105009; adjusted 0.6676099
- wealthindex: unadjusted 0.702451; adjusted 0.6450516

## Random effect model

```
# Simulation parameters
nHH <- 1000      # Number of HH (primary units) per year

nyrs <- 10      # Number of years to simulate
yrs <- 2000 + c(1:nyrs) # Years to simulate
N <- nyrs * nHH

## HH random effect sd
hhSD <- 0.5

# Generate dataset template
```

```

temp_df <- (data.frame(hhid = rep(c(1:nHH), each = nyrs)
  , years = rep(yrs, nHH)
  , age = runif(n=N, age_min, age_max)
  , wealthindex = rnorm(n = N, 0, 1)
)
)

# Simulate HH-level random effects (residual error)
hhRE <- rnorm(nHH, hhSD)
temp_df$hhRE <- hhRE[temp_df$hhid]

sim_df <- (temp_df
  %>% mutate(eta = beta0 + betaA * age + betaW * wealthindex + hhRE
    , status = rbinom(N, 1, plogis(eta))
  )
  %>% select(-eta)
)
true_prop_reff <- mean(sim_df$status)
print(true_prop)

## [1] 0.70092
print(head(sim_df, 50))

```

##	hhid	years	age	wealthindex	hhRE	status
## 1	1	2001	0.7639604	-0.344855381	1.7135445	1
## 2	1	2002	0.8449101	0.003245716	1.7135445	1
## 3	1	2003	0.3191792	-0.400368171	1.7135445	1
## 4	1	2004	0.5275991	0.493046821	1.7135445	1
## 5	1	2005	0.4625588	-0.825271859	1.7135445	1
## 6	1	2006	0.2967984	-0.230758439	1.7135445	1
## 7	1	2007	0.6244001	-0.216376416	1.7135445	0
## 8	1	2008	0.3798349	-0.950330918	1.7135445	1
## 9	1	2009	0.6987809	1.670674922	1.7135445	1
## 10	1	2010	0.9802189	-1.275773471	1.7135445	0
## 11	2	2001	0.2410522	0.523871564	0.2985271	1
## 12	2	2002	0.8135433	1.117013232	0.2985271	1
## 13	2	2003	0.6907868	-1.002174686	0.2985271	1
## 14	2	2004	0.7820465	0.424314274	0.2985271	1
## 15	2	2005	0.2713455	-0.722798175	0.2985271	0
## 16	2	2006	0.7629002	-0.143273633	0.2985271	0
## 17	2	2007	0.6066807	0.436517100	0.2985271	1
## 18	2	2008	0.7122677	-0.646761524	0.2985271	0
## 19	2	2009	0.7586238	1.005114038	0.2985271	1
## 20	2	2010	0.4696731	-0.499505689	0.2985271	1
## 21	3	2001	0.6544598	1.360953916	0.6099605	1
## 22	3	2002	0.3199612	0.559676147	0.6099605	0
## 23	3	2003	0.7047462	0.313300235	0.6099605	0
## 24	3	2004	0.9865054	0.335200959	0.6099605	1
## 25	3	2005	0.4032590	-1.233651878	0.6099605	1
## 26	3	2006	0.6469325	-1.444368932	0.6099605	1
## 27	3	2007	0.2042122	-1.068794572	0.6099605	1
## 28	3	2008	0.7705646	2.199997885	0.6099605	1
## 29	3	2009	0.5240763	-1.464230742	0.6099605	0

```
## 30    3  2010 0.6841513  0.259446519  0.6099605    1
## 31    4  2001 0.5142064 -0.570075632  0.6524803    1
## 32    4  2002 0.4027693  0.861234446  0.6524803    0
## 33    4  2003 0.8963569 -0.632076733  0.6524803    1
## 34    4  2004 0.5185198 -0.961957492  0.6524803    1
## 35    4  2005 0.4009540  1.494735570  0.6524803    1
## 36    4  2006 0.5912513 -1.558898248  0.6524803    0
## 37    4  2007 0.5228832  0.911760283  0.6524803    1
## 38    4  2008 0.2940842 -1.166798623  0.6524803    1
## 39    4  2009 0.9178963 -0.513016059  0.6524803    1
## 40    4  2010 0.7462044 -1.786672955  0.6524803    1
## 41    5  2001 0.8687909 -0.835155399 -0.2066059    1
## 42    5  2002 0.9870285  0.705113579 -0.2066059    1
## 43    5  2003 0.4755328  0.044927990 -0.2066059    1
## 44    5  2004 0.7545575 -0.010572632 -0.2066059    1
## 45    5  2005 0.7341992 -0.801877636 -0.2066059    1
## 46    5  2006 0.4047338  0.715316850 -0.2066059    0
## 47    5  2007 0.6860020 -1.557316158 -0.2066059    0
## 48    5  2008 0.8667148  0.833676669 -0.2066059    1
## 49    5  2009 0.5129266  0.881191045 -0.2066059    0
## 50    5  2010 0.9262759  0.456788948 -0.2066059    1
```

## Fit model

```
reff_mod <- glmmTMB(status ~ age + wealthindex + (1|hhid)
  , data = sim_df
  , family = binomial(link = "logit")
)
```

```
## Tidy coef estimates
reff_coef_df <- (broom.mixed::tidy(reff_mod, conf.int=TRUE)
  %>% mutate(term = gsub("\\(|\\)", "", term))
  %>% filter(effect=="fixed")
)
```

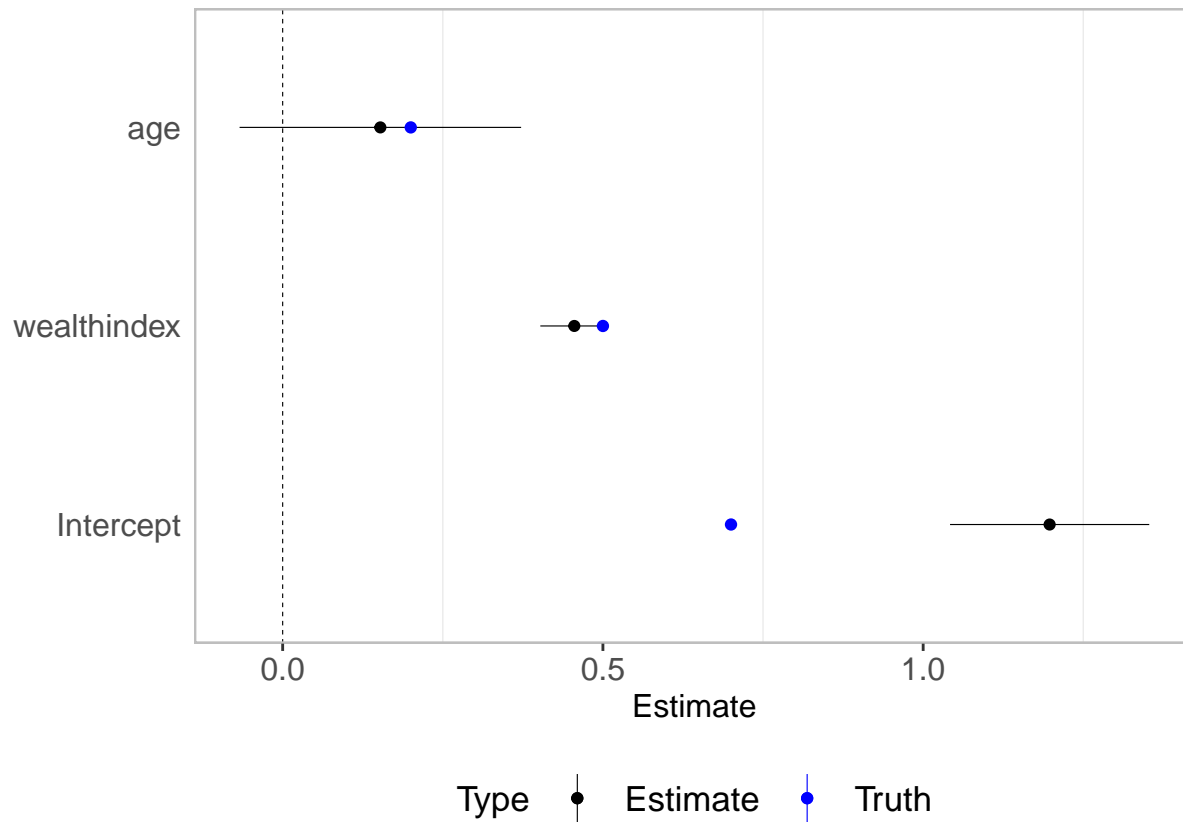
```
## Registered S3 method overwritten by 'broom.mixed':
##   method      from
##   tidy.gamlss broom
```

```
print(reff_coef_df)
```

```
## # A tibble: 3 x 10
##   effect component group term      estimate std.error statistic  p.value conf.low
##   <chr>   <chr>      <chr> <chr>      <dbl>     <dbl>     <dbl>    <dbl>    <dbl>
## 1 fixed   cond        <NA> Interce~    1.20     0.0793     15.1  1.84e-51    1.04
## 2 fixed   cond        <NA> age         0.153     0.112      1.36  1.74e- 1   -0.0673
## 3 fixed   cond        <NA> wealthi~    0.455     0.0271     16.8  1.87e-63    0.402
## # ... with 1 more variable: conf.high <dbl>
```

```
reff_coef_plot <- (plotEsize(reff_coef_df)
  + geom_point(data=true_beta_df, aes(x=term, y=estimate, colour="Truth"))
  + labs(colour="Type")
  + scale_colour_manual(values=c("black", "blue"))
)
print(reff_coef_plot)
```





Variable effect plots

- Age

```
## varpred way
reff_vareff_age <- varpred(reff_mod, "age", isolate=FALSE, modelname="varpred")

## Pop. average
reff_vareff_age_pop <- varpred(reff_mod, "age", isolate=TRUE, pop.ave=TRUE, modelname = "Pop. ave")

## Bias adjust
reff_vareff_age_adjust <- varpred(reff_mod, "age", isolate=TRUE, bias.adjust=TRUE, modelname = "Bias adj.")

vareff_age <- reff_vareff_age
vareff_age$preds <- do.call("rbind", list(vareff_age$preds, reff_vareff_age_pop$preds, reff_vareff_age_adjust$preds))
age_plot <- (plot(vareff_age)
  + labs(y="Prob. of improved \n service", colour="Model")
  + geom_hline(yintercept=true_prop_reff, lty=2, colour="grey")
  + scale_colour_manual(values=c("black", "blue", "red"))
  + theme(legend.position="bottom")
)

## Scale for 'colour' is already present. Adding another scale for 'colour',
## which will replace the existing scale.
```

- Wealth index

```
# Wealth index
## varpred
reff_vareff_wealthindex <- varpred(reff_mod, "wealthindex", isolate=TRUE, modelname="varpred")
```

```

## Pop. average
reff_vareff_wealthindex_pop <- varpred(reff_mod, "wealthindex", isolate=TRUE, pop.ave=TRUE, modelname="")

## Bias adjust
reff_vareff_wealthindex_adjust <- varpred(reff_mod, "wealthindex", isolate=TRUE, bias.adjust=TRUE, modelname="")

vareff_wealthindex <- reff_vareff_wealthindex
vareff_wealthindex$preds <- do.call("rbind", list(vareff_wealthindex$preds, reff_vareff_wealthindex_pop$preds))
wealthindex_plot <- (plot(vareff_wealthindex)
  + labs(y="", colour="Model")
  + geom_hline(yintercept=true_prop_reff, lty=2, colour="grey")
  + scale_colour_manual(values=c("black", "blue", "red"))
  + theme(legend.position="bottom")
)

## Scale for 'colour' is already present. Adding another scale for 'colour',
## which will replace the existing scale.
ggarrange(age_plot, wealthindex_plot, common.legend=TRUE)

```

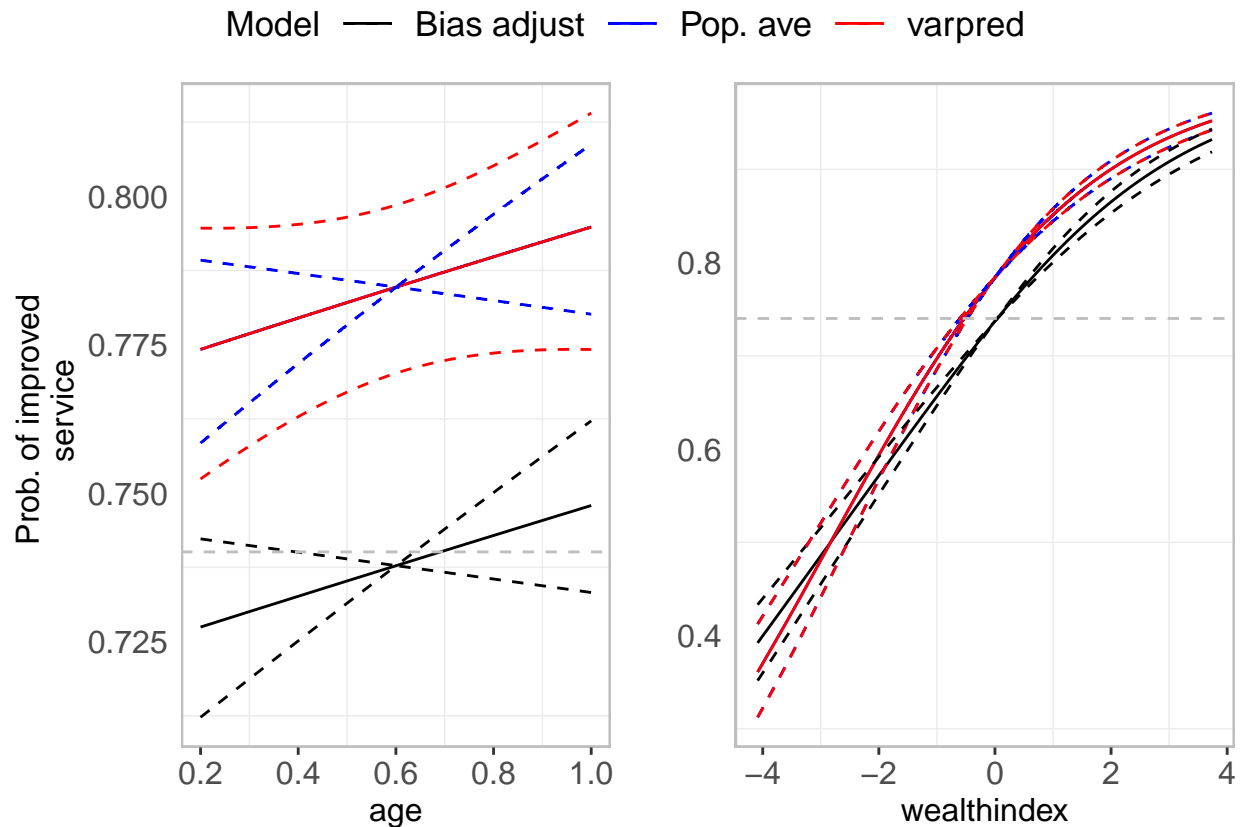


Figure 2: A comparison of population averaged, varpred-based and bias-adjusted predictions. For each of the predictors, the population averaging and varpred gives similar estimates, and slightly off the truth. However, when we apply bias-adjustment, the estimates are very precise (i.e., close to the truth).