

Continuous Simulation

Tucker Morgan - tlm2152

2/16/2022

Generating Covariates and Finding Coefficients in Propensity Model

```
set.seed(100)
covariate_coef <- function(desired_prop, cov_df) {

  alpha_0 = log(desired_prop/(1 - desired_prop))

  coef_L1 <- sample(cov_df$L1, 10000, replace = TRUE)
  coef_L2 <- sample(cov_df$L2, 10000, replace = TRUE)
  # coef_L3 <- sample(cov_df$L3, 10000, replace = TRUE)

  A_logit <- vector(mode = "list", length = length(coef_L1))
  p_A <- vector(mode = "numeric", length = length(coef_L1))

  u <- alpha_0 + coef_L1[1]*cov_df$L1 + coef_L2[1]*cov_df$L2 #- coef_L3[1]*cov_df$L3

  p_A[1] <- mean(exp(u)/(1 + exp(u)))

  tol <- 0.001

  i = 1

  while (abs(p_A[i] - 0.14) > tol) {

    i = i + 1

    A_logit[[i]] <-
      alpha_0 + coef_L1[i]*cov_df$L1 + coef_L2[i]*cov_df$L2 #- coef_L3[i]*cov_df$L3

    p_A[i] <- mean(exp(A_logit[[i]])/(1 + exp(A_logit[[i]])))

    if (abs(p_A[i] - 0.14) < tol) {
      mean_treated_proportion <- p_A[i]
      desired_coef_L1 <- coef_L1[i]
      desired_coef_L2 <- coef_L2[i]
      # desired_coef_L3 <- coef_L3[i]
    }
  }
}
```

```

    if (i > length(coef_L1)) {
      stop("You need better coverage, mate.")
    }
  }

  return(tibble(alpha_0, mean_treated_proportion,
    desired_coef_L1, desired_coef_L2))#, desired_coef_L3))
}

```

Generating 100 Sub-Populations

```

seed_vec <- rnorm(100000, mean = 0, sd = 100) %>% round(0) %>% unique()

generate_no_boot_data <- function(n, size = 5000, seeds = seed_vec) {

  df <- list()

  cov_df <- list()

  pb <- progress_bar$new(format = "generating data... [:bar]", total = n)

  for (i in 1:n) {
    pb$tick()
    set.seed(seeds[i])

    set.seed(seeds[i])
    pre_data <- defData(varname = "L1", formula = "0", variance = 1,
      dist = "normal")

    pre_data <- defData(pre_data, varname = "L2", formula = "0", variance = 10,
      dist = "normal")

    pre_data <- defData(pre_data, varname = "L3", formula = "0", variance = 10,
      dist = "normal")

    cov_df[[i]] <- genData(5000, pre_data)

    cov_coef_df <- covariate_coef(0.14, cov_df[[i]])

    L1_coef <- cov_coef_df$desired_coef_L1
    L2_coef <- cov_coef_df$desired_coef_L2
    #L3_coef <- cov_coef_df$desired_coef_L3

    pre_data <- defData(pre_data, varname = "L1_coef", formula = L1_coef)
    pre_data <- defData(pre_data, varname = "L2_coef", formula = L2_coef)
    #pre_data <- defData(pre_data, varname = "L3_coef", formula = L3_coef)

    pre_data <- defData(pre_data, varname = "A",
      formula = "-1.815 + L1_coef*L1 + L2_coef*L2", # + L3_coef*L3",
      dist = "binary", link = "logit")
  }
}

```

```

pre_data <- defData(pre_data, varname = "Y",
                    formula = ".5 + 0.15*A + 2*L2 - 1*L3" ,
                    dist = "nonrandom")

df[[i]] <- genData(size, pre_data)
df[[i]] <- df[[i]] %>% select(-L1_coef, -L2_coef)#, -L3_coef)
}
return(df)
}

ate_true = 0.15
no_boot_list <- generate_no_boot_data(100)

```

Let's take a look at one of our "no-boot" aka sub-population data sets.

```

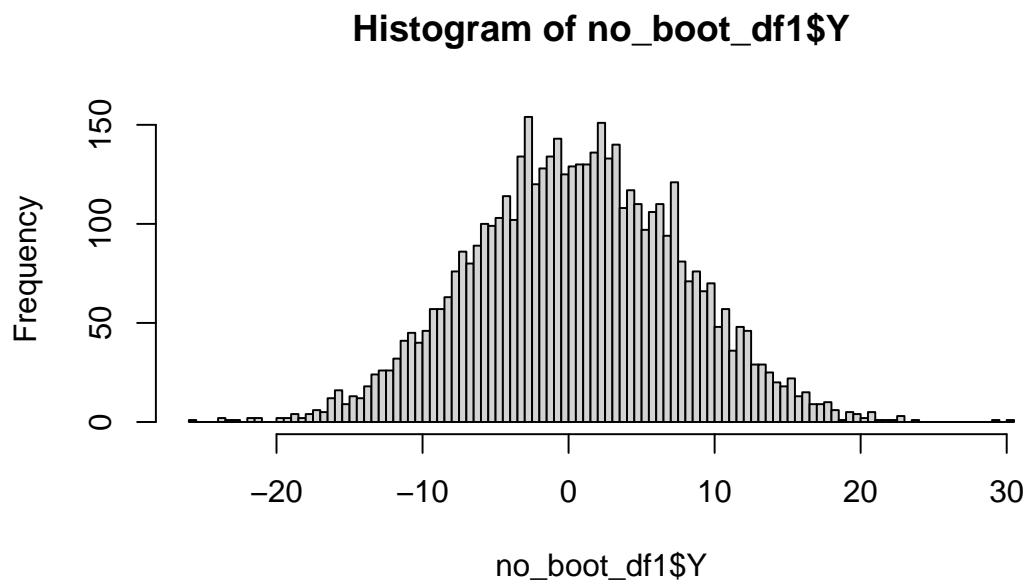
no_boot_df1 <- no_boot_list[[1]]

sum(no_boot_df1$A) / nrow(no_boot_df1) # similar to 0.14

```

```
## [1] 0.1452
```

```
hist(no_boot_df1$Y, breaks = 100) # continuous distribution of outcome
```



Implementing Nearest-Neighbor Matching

```
df <- no_boot_list
```

```

# could possible turn this into a function later.
matched_df <- list()

for (i in 1:length(df)) {

  matched <- matchit(A ~ L1 + L2 + L3,
                    data = df[[i]],
                    distance = "glm",
                    link = "logit",
                    method = "nearest",
                    ratio = 1) # perform NNM

  matched_df[[i]] <- match.data(matched, distance = "ps")
}

```

Again, we'll look at one of our matched sub-pops.

```

matched_df1 <- matched_df[[1]]
summary(matched_df1)

```

```

##           id           L1           L2           L3
## Min.      : 6    Min.    :-3.2858    Min.    :-10.8706    Min.    :-11.94392
## 1st Qu.:1257    1st Qu.: -0.5480    1st Qu.: -2.0860    1st Qu.: -2.24535
## Median :2520    Median : 0.1383    Median : 0.1545    Median : -0.03596
## Mean    :2506    Mean     : 0.1215    Mean     : 0.1254    Mean     : -0.07649
## 3rd Qu.:3739    3rd Qu.: 0.7990    3rd Qu.: 2.2444    3rd Qu.: 2.14221
## Max.    :5000    Max.     : 3.8761    Max.     : 10.8396    Max.     : 10.87334
##
##           A           Y           ps           weights           subclass
## Min.      :0.0    Min.    :-23.2913    Min.     :0.09977    Min.     :1    1      : 2
## 1st Qu.:0.0    1st Qu.: -4.0356    1st Qu.:0.13487    1st Qu.:1    2      : 2
## Median :0.5    Median : 0.8846    Median :0.14679    Median :1    3      : 2
## Mean    :0.5    Mean     : 0.9023    Mean     :0.14802    Mean     :1    4      : 2
## 3rd Qu.:1.0    3rd Qu.: 5.7835    3rd Qu.:0.16062    3rd Qu.:1    5      : 2
## Max.    :1.0    Max.     : 29.3141    Max.     :0.25479    Max.     :1    6      : 2
##                                     (Other):1440

```

```

str(matched_df1)

```

```

## Classes 'matchdata', 'data.table' and 'data.frame': 1452 obs. of 9 variables:
## $ id      : int 6 12 13 15 26 29 30 34 47 50 ...
## $ L1      : num 1.658 0.861 -0.198 -0.597 -0.866 ...
## $ L2      : num 4.73 3.35 3.78 6.02 4.54 ...
## $ L3      : num -6.843 -3.351 4.76 0.768 -2.605 ...
## $ A       : int 0 1 0 0 0 0 1 1 1 1 ...
## $ Y       : num 16.8 10.7 3.3 11.8 12.2 ...
## $ ps      : num 0.2 0.175 0.151 0.155 0.147 ...
## $ weights : num 1 1 1 1 1 1 1 1 1 1 ...
## $ subclass: Factor w/ 726 levels "1","2","3","4",...: 204 31 395 251 133 316 322 384 593 640 ...
## - attr(*, ".internal.selfref")=<externalptr>
## - attr(*, "distance")= chr "ps"
## - attr(*, "weights")= chr "weights"
## - attr(*, "subclass")= chr "subclass"

```

The Simple Bootstrap

```
# creating the tibble to apply map function
matched_tib <-
  tibble(data = matched_df)

# ### function to iterate glm over a list, to be used in purr:map ###
# returns tibble of parameter estimates and standard errors.

outcome_model_list <- function(list) {
  tib_coef <- tibble()
  pb3$tick()
  for (i in 1:length(list)) {
    mod <- glm(Y ~ A + ps,
              data = list[[i]],
              weights = weights) %>% summary()
    coefs <- mod$coefficients[2,1:2]
    tib_coef <- bind_rows(tib_coef, tibble(estimate = coefs[1], se = coefs[2]))
  }
  return(tib_coef)
}

# ### input matched dataframe, output however many bootstrapped samples you want ###
# first, set seed vector for reproducibility

# now, define function

seed_vec_2 <- rnorm(100000, mean = 0, sd = 10000) %>% round(0) %>% unique()

simple_boot <- function(df, n, size = 500, seeds = seed_vec_2){
  boots <- list()
  pb2$tick()
  for (i in 1:n) {
    set.seed(seeds[i])
    boots[[i]] <-
      df %>%
      filter(subclass %in% sample(levels(subclass),
                                size,
                                replace = TRUE))
  }
  return(boots)
}

# adding progress bars for sanity
pb2 <- progress_bar$new(format = "bootstrapping... [:bar]", total = nrow(matched_tib))
pb3 <- progress_bar$new(format = "performing glm... [:bar]", total = nrow(matched_tib))

# creating booted tibbles, applying functions through purr:map.
boot_tib <-
  matched_tib %>%
  mutate(
```

```

  boots = map(.x = data, ~simple_boot(.x, n = 1000))
) %>%
mutate(coef = map(.x = boots, ~outcome_model_list(.x)))

```

```

boot_estimates <-
  boot_tib %>%
  mutate(seq = seq(1:nrow(boot_tib))) %>%
  select(coef, seq) %>% unnest(coef)

```

Summary of 1000 Bootstraps in 100 Sub-Populations

```

boot_result <-
  boot_estimates %>%
  group_by(seq) %>%
  summarize(avg_trt_eff = mean(estimate), sd_ate = sd(estimate))

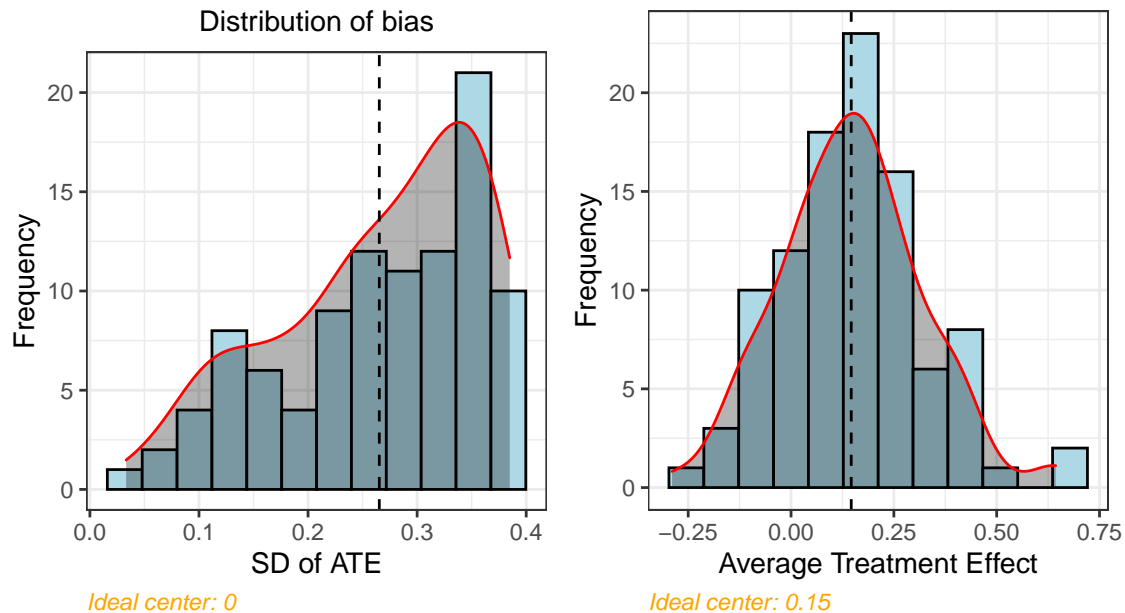
fig1 <-
  boot_result %>%
  ggplot(aes(x = sd_ate, color = sd_ate)) +
  geom_histogram(fill = "light blue", bins = 12, color = "black") +
  geom_density(aes(y = ..density..*4), colour = "red",
    fill = "black", alpha = 0.3) +
  geom_vline(xintercept = mean(boot_result$sd_ate), linetype = "dashed") +
  labs(title = "SD of ATE from 1000 Bootstraps in 100 Sub-Populations",
    subtitle = "Distribution of bias",
    caption = "Ideal center: 0", x = "SD of ATE", y = "Frequency") +
  theme(
    plot.title = element_text(color = "blue", size = 11, face = "bold"),
    plot.subtitle = element_text(color = "black"),
    plot.caption = element_text(color = "orange", face = "italic")
  )

fig2 <-
  boot_result %>%
  ggplot(aes(x = avg_trt_eff)) +
  geom_histogram(fill = "light blue", bins = 12, color = "black") +
  geom_density(aes(y = ..density..*8), colour = "red",
    fill = "black", alpha = 0.3) +
  geom_vline(xintercept = mean(boot_result$avg_trt_eff), linetype = "dashed") +
  labs(title = "Distribution of ATE in 1000 Bootstraps of 100 Sub-Populations",
    caption = "Ideal center: 0.15", x = "Average Treatment Effect", y = "Frequency") +
  theme(
    plot.title = element_text(color = "blue", size = 11, face = "bold"),
    plot.caption = element_text(color = "orange", face = "italic")
  )

plot_grid(fig1, fig2)

```

Distribution of ATE from 1000 Bootstraps in 1000 Subsets vs Distribution of ATE in 1000 Bootstraps of 100 Subsets



Confidence Intervals Coverage Rates

```
cvg_rate <- function(df){
  res = df %>%
    mutate(ci_low = avg_trt_eff - 1.96*sd_ate,
           ci_high = avg_trt_eff + 1.96*sd_ate,
           covered = case_when(
             ci_low <= ate_true & ci_high >= ate_true ~ 1,
             TRUE ~ 0
           ))

  return(sum(res$covered) / nrow(res))
}

cvg_plot <- function(df){
  res = df %>%
    mutate(ci_low = avg_trt_eff - 1.96*sd_ate,
           ci_high = avg_trt_eff + 1.96*sd_ate,
           covered = case_when(
             ci_low <= ate_true & ci_high >= ate_true ~ 1,
             TRUE ~ 0
           ))

  plot = res %>%
    ggplot(aes(x = avg_trt_eff, y = seq)) +
    geom_point() +
    geom_errorbar(aes(xmin = ci_low, xmax = ci_high)) +
    geom_vline(xintercept = ate_true, linetype = "dashed")

  return(plot)
}
```

```
}  
  
cvg_rate(boot_result)
```

```
## [1] 0.99
```

```
cvg_plot(boot_result)
```

