

# Session 3

## Sample Size Estimation

---

Matt Denwood

2024-01-30

## A note on coding styles

My code looks a bit different to Giles's code, e.g.:

```
# If necessary:  
## install.packages(c("tidyverse", "pbapply"))  
  
library("tidyverse")  
library("pbapply")
```

## A note on coding styles

My code looks a bit different to Giles's code, e.g.:

```
# If necessary:  
## install.packages(c("tidyverse", "pbapply"))  
  
library("tidyverse")  
library("pbapply")
```

REMEMBER: the coding style is not important as long as the output is the same!

## **Background to sample size calculations**

---

# Power calculation

Power is defined as the proportion of experiments that can be expected to give p-values of  $\leq 0.05$  (or whatever alpha is chosen), conditional on the specified parameters.

Power calculations can be done using:

## 1. Approximation methods, e.g. power.t.test:

```
power.t.test(n = 150, delta = 0.25, sd = 1)
##
##      Two-sample t test power calculation
##
##              n = 150
##            delta = 0.25
##              sd = 1
##      sig.level = 0.05
##            power = 0.5785239
##      alternative = two.sided
##
## NOTE: n is number in *each* group
```

## 2. Numerical methods i.e. by simulation:

A function to simulate data, then calculate and return a p-value:

```
p_fun <- function(parameters){  
  stopifnot(is.data.frame(parameters), nrow(parameters)==1L, "Size" %in% names(parameters))  
  sample1 <- rnorm(parameters$Size, mean=0, sd=1)  
  sample2 <- rnorm(parameters$Size, mean=0.25, sd=1)  
  parameters |>  
    mutate(P_val = t.test(sample1, sample2)$p.value)  
}
```

```
p_fun(tibble(Size = 150L))  
## # A tibble: 1 x 2  
##   Size P_val  
##   <int> <dbl>  
## 1    150 0.0171
```

There is randomness so this will be different every time it is run:

```
p_fun(tibble(Size = 150L))  
## # A tibble: 1 x 2  
##   Size P_val  
##   <int> <dbl>  
## 1   150 0.0109  
p_fun(tibble(Size = 150L))  
## # A tibble: 1 x 2  
##   Size P_val  
##   <int> <dbl>  
## 1   150 0.0119
```

So we must run it several times (e.g. 1000):

```
tibble(Iteration = seq_len(1000L), Size = 150L) |>  
  group_split(Iteration, Size) |>  
  lapply(p_fun) |>  
  bind_rows() ->  
  pvals
```



So we must run it several times (e.g. 1000):

```
tibble(Iteration = seq_len(1000L), Size = 150L) |>
  group_split(Iteration, Size) |>
  lapply(p_fun) |>
  bind_rows() ->
  pvals
```

And we calculate the power like so:

```
pvals |>
  group_by(Size) |>
  summarise(Power = sum(P_val <= 0.05) / n(), .groups="drop")
## # A tibble: 1 x 2
##   Size Power
##   <int> <dbl>
## 1    150 0.581
```

# Sample size estimation

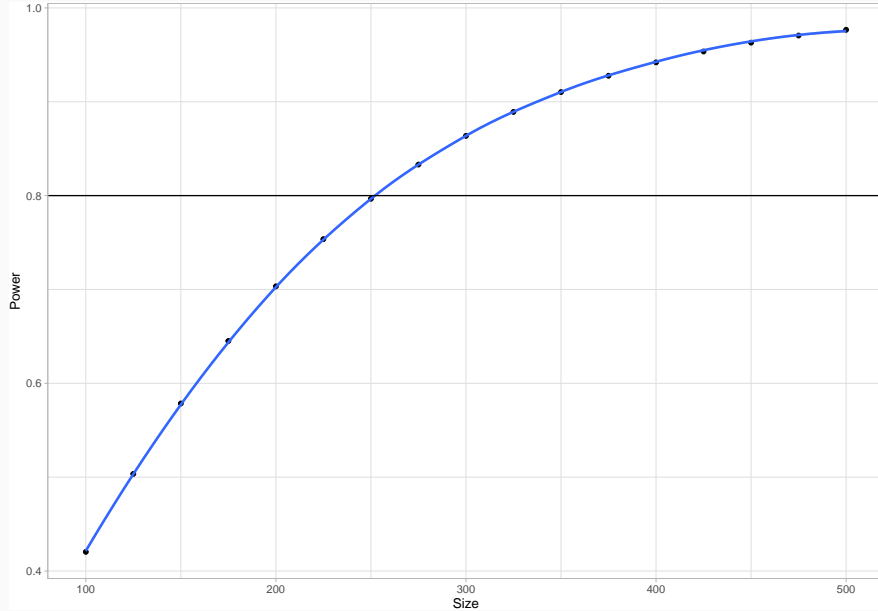
The goal is typically to find the minimum sample size that corresponds to  $\geq 80\%$  power, for a specified set of parameters. This can be done in one of two ways:

1. Using approximation methods directly i.e.:

```
power.t.test(n = NULL, delta = 0.25, sd = 1, power = 0.8)
##
##      Two-sample t test power calculation
##
##              n = 252.1281
##            delta = 0.25
##              sd = 1
##      sig.level = 0.05
##            power = 0.8
## alternative = two.sided
##
## NOTE: n is number in *each* group
```

## 2. By trying different sample sizes (using either approximation methods or simulation):

```
tibble(Size = seq(100, 500, by=25)) |>
  group_split(Size) |>
  lapply(function(parameters){
    parameters |>
      mutate(Power = power.t.test(n = parameters$Size, delta = 0.25, sd = 1)$power)
  }) |>
  bind_rows() ->
  power_estimates
```



## **Sample size calculation for LCM**

---

## Determining the objective

Let's take a simple 2-test, 2-population Hui-Walter model as an example.

- There is (usually) no hypothesis to test, so we don't have a concept of power
- Our objective is to have narrow posterior (95%) credible intervals (“precision”\*)
- But for which parameter?
- Note: this is not exactly the same as the usual statistical definition of precision

# Types of parameter

In general:

- Experimental parameter: structural things that we can control
- Parameter of interest: things we want to estimate
- Nuisance parameters: everything else

# Types of parameter

In general:

- Experimental parameter: structural things that we can control
- Parameter of interest: things we want to estimate
- Nuisance parameters: everything else

For t-tests:

- Experimental parameter: sample size
- Parameter of interest: difference in means
- Nuisance parameters: standard deviation



In general:

- Experimental parameter: structural things that we can control
- Parameter of interest: things we want to estimate
- Nuisance parameters: everything else

For Hui-Walter models:

- Experimental parameters: number of samples from each population (and maybe number of populations)
- Parameters of interest: sensitivities, specificities
- Nuisance parameters: prevalences (and maybe correlation terms for  $>2$  tests)

## Determining the objective

We have 4 parameters of interest:  $2 \times \text{Se}$ ,  $2 \times \text{Sp}$ . Which are most important?

- Is sensitivity the sole objective? NB: high-prevalence population will be prioritised.
- Is specificity the sole objective? NB: low-prevalence population will be prioritised.
- Are sensitivity and specificity equally important? If so, do we average the relative or absolute size of their 95% CI?

## Determining the objective

We have 4 parameters of interest:  $2 \times \text{Se}$ ,  $2 \times \text{Sp}$ . Which are most important?

- Is sensitivity the sole objective? NB: high-prevalence population will be prioritised.
- Is specificity the sole objective? NB: low-prevalence population will be prioritised.
- Are sensitivity and specificity equally important? If so, do we average the relative or absolute size of their 95% CI?
- Would Youden's index be easier?

## Determining the objective

We have 4 parameters of interest:  $2 \times \text{Se}$ ,  $2 \times \text{Sp}$ . Which are most important?

- Is sensitivity the sole objective? NB: high-prevalence population will be prioritised.
- Is specificity the sole objective? NB: low-prevalence population will be prioritised.
- Are sensitivity and specificity equally important? If so, do we average the relative or absolute size of their 95% CI?
- Would Youden's index be easier?
- One or both tests?

## Determining the objective

We have 4 parameters of interest:  $2 \times \text{Se}$ ,  $2 \times \text{Sp}$ . Which are most important?

- Is sensitivity the sole objective? NB: high-prevalence population will be prioritised.
- Is specificity the sole objective? NB: low-prevalence population will be prioritised.
- Are sensitivity and specificity equally important? If so, do we average the relative or absolute size of their 95% CI?
- Would Youden's index be easier?
- One or both tests?

My suggestion: define “precision” as the average width of 95% CI for Youden's index (for either one or both tests)

# Simulating data

Best as a function that takes population-level and test-level inputs:

```
simulation_fun <- function(populations, tests){  
  stopifnot(  
    is.data.frame(populations),  
    nrow(populations) >= 2L,  
    c("N", "Prev") %in% names(populations),  
    populations$N >= 1L,  
    populations$N %% 1 == 0,  
    populations$Prev >= 0, populations$Prev <= 1  
  )  
  stopifnot(  
    is.data.frame(tests),  
    nrow(tests) >= 2L,  
    c("Se", "Sp") %in% names(tests),  
    tests$Se >= 0, tests$Se <= 1,  
    tests$Sp >= 0, tests$Sp <= 1  
  )  
  ## Do some stuff like from session 2 and return the simulated dataset  
  ## See the exercise for a complete function  
}
```

Output looks like:

```
tests <- tribble(
  ~Se, ~Sp,
  0.8, 0.99,
  0.9, 0.95
)

populations <- tribble(
  ~N, ~Prev,
  100, 0.1,
  100, 0.4
)

(data <- simulation_fun(populations, tests))
##      [,1] [,2]
## [1,]   89   59
## [2,]    0    4
## [3,]    3    6
## [4,]    8   31
```

# Analyzing data

Also best as a function taking the data, burnin and sample iterations as inputs:

```
analysis_fun <- function(data, burnin=1000L, sample=5000L){  
  stopifnot(is.matrix(data), nrow(data)==4L, ncol(data)==2L, data>=0L)  
  ## Do some stuff like from session 2 to analyse the data  
  ## See the exercise for a complete function  
  results |>  
    summary(vars=c("youden","se","sp")) |>  
    as.data.frame() |>  
    rownames_to_column("Variable")  
}
```



## Output looks like:

```
analysis_fun(data)
## Loading required namespace: rjags
##      Variable   Lower95     Median   Upper95      Mean
## 1 youden[1] 0.7342350 0.8628030 0.9906417 0.8580342
## 2 youden[2] 0.7716093 0.8805874 0.9761025 0.8754820
## 3      se[1] 0.7518793 0.8753646 0.9998213 0.8705460
## 4      se[2] 0.8179392 0.9128422 0.9999820 0.9073788
## 5      sp[1] 0.9649287 0.9906899 0.9999998 0.9874882
## 6      sp[2] 0.9296196 0.9708700 0.9999972 0.9681032
##           SD      Mode      MCerr MC%ofSD SSeff
## 1 0.06855032 0.8718323 0.0009487945      1.4 5220
## 2 0.05390644 0.8874831 0.0006481892      1.2 6916
## 3 0.06760940 0.8793672 0.0009359165      1.4 5218
## 4 0.04945453 0.9231452 0.0005719271      1.2 7477
## 5 0.01130297 0.9958692 0.0001682663      1.5 4512
## 6 0.02087240 0.9780313 0.0003053392      1.5 4673
##           AC.10      psrf
## 1 0.018980356 0.9999562
## 2 0.003457944 1.0001022
## 3 0.020654609 0.9999569
## 4 0.009074103 1.0003880
## 5 0.004984732 0.9999591
## 6 0.026690476 1.0001098
```

## A quick note on label switching

I now recommend this method of specifying minimally informative priors:

```
model{  
  ### Rest of the model as usual  
  
  for(t in 1:2){  
    se[t] ~ dbeta(2,1)  
    sp[t] ~ dbeta(2,1)  
    youden[t] <- se[t]+sp[t]-1.0  
    AcceptTest[t] ~ dbern(ifelse(youden[t] >= 0.0, 1, 0))  
  }  
  #data# AcceptTest  
}  
  
AcceptTest <- c(1,1)
```

## A quick note on label switching

I now recommend this method of specifying minimally informative priors:

```
model{  
  ### Rest of the model as usual  
  
  for(t in 1:2){  
    se[t] ~ dbeta(2,1)  
    sp[t] ~ dbeta(2,1)  
    youden[t] <- se[t]+sp[t]-1.0  
    AcceptTest[t] ~ dbern(ifelse(youden[t] >= 0.0, 1, 0))  
  }  
  #data# AcceptTest  
}  
  
AcceptTest <- c(1,1)
```

More on this during my presentation tomorrow!

## Combining the two

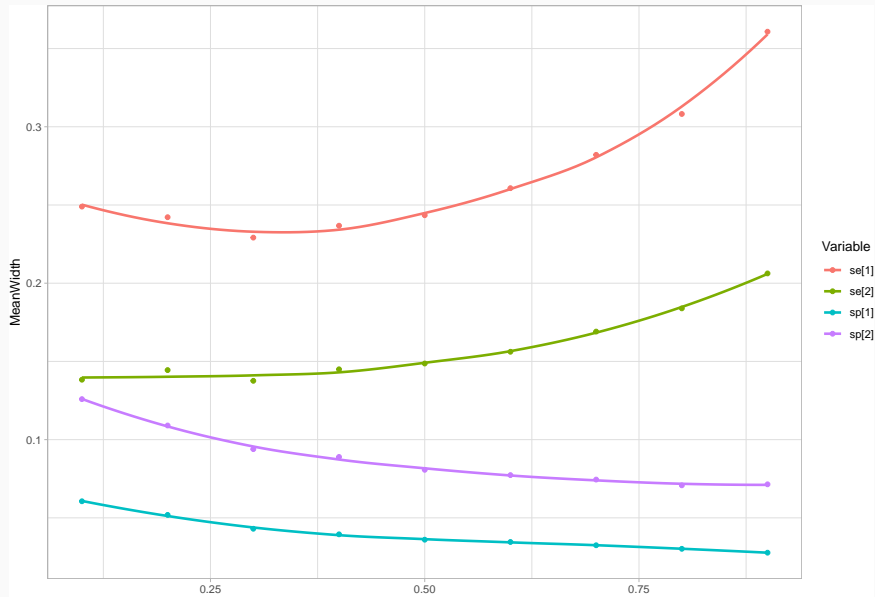
```
summary_fun <- function(parameters, iterations, cl=NULL, burnin=1000L, sample=5000L){
  stopifnot(is_tibble(parameters))
  stopifnot(iterations >= 1L)

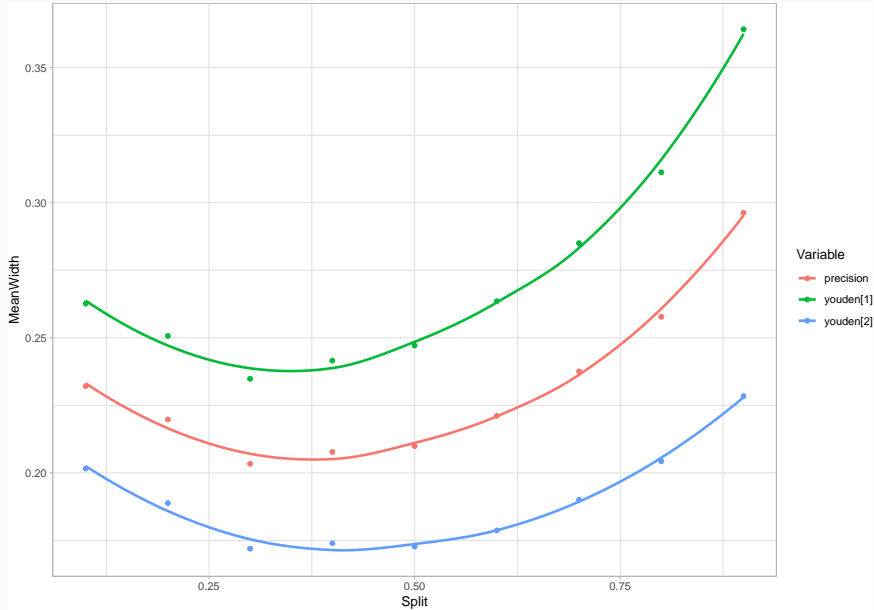
  parameters |>
    mutate(ParameterSet = row_number()) |>
    expand_grid(Iteration = seq_len(iterations)) |>
    rowwise() |>
    group_split() |>
    pblapply(function(x){
      simulation_fun(x$Populations[[1]], x$Tests[[1]]) |>
        analysis_fun(burnin=burnin, sample=sample) |>
        bind_cols(x)
    }, cl=cl) |>
    bind_rows() |>
    mutate(WidthCI = Upper95 - Lower95) |>
    group_by(ParameterSet, Variable) |>
    summarise(MeanEst = mean(Mean), MeanWidth = mean(WidthCI), MeanLCI = mean(Lower95),
      ↪ MeanUCI = mean(Upper95), .groups="drop") |>
    full_join(parameters |> mutate(ParameterSet = row_number()), by="ParameterSet")
}
```

Output looks like:

```
parameters <- tibble(Populations=list(populations), Tests=list(tests))
summary_fun(parameters, 10L)
## # A tibble: 6 x 8
##   ParameterSet Variable  MeanEst MeanWidth MeanLCI MeanUCI
##         <int> <chr>      <dbl>     <dbl>   <dbl>   <dbl>
## 1             1 se[1]      0.833     0.293   0.687   0.980
## 2             1 se[2]      0.909     0.183   0.812   0.994
## 3             1 sp[1]      0.979     0.0510  0.949   1.00
## 4             1 sp[2]      0.942     0.107   0.888   0.994
## 5             1 youden[1]  0.812     0.305   0.659   0.964
## 6             1 youden[2]  0.851     0.226   0.734   0.960
## # i 2 more variables: Populations <list>, Tests <list>
```

# Visualising the results





It is easier to see a balance between these using Youden's index or our “precision”.

## Exercises

1. Look at the function code given in the HTML file: read it through and make sure you understand it!
2. Re-create the “visualising the results” plot I have given above. Test #1 has Se/Sp of 0.8/0.99 and Test #2 has Se/Sp of 0.9/0.95. Prevalences in the two populations are 10% and 40%, and the total sample size is 500.
3. Assuming that the optimum distribution of tests between the two populations is always around 40% / 60% for these test/prevalence estimates, create a plot showing how precision increases with total sample size of between 100 and 1000.
4. How much does the precision also depend on the parameters of interest (diagnostic test performance) and nuisance parameters (prevalences)?