

A General Approach to Sample Size Analysis

Mihai A. Constantin
Noémi N. K. Schuurman
Jeroen K. Vermunt

mihai@mihaiconstantin.com

So Far

We've seen:

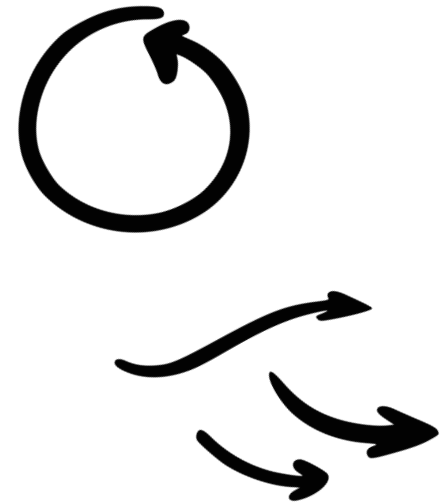
- what **sample size planning** is and why it matters
- two **criteria** for searching for an **optimal sample size**
 - statistical power
 - predictive accuracy
- two **approaches** for conducting **sample size analysis**
 - analytical
 - simulation
- applications to time series models, i.e., $AR(1)$ and $VAR(1)$



So Next

We'll talk about:

- the **requirements** of **simulation**-based **sample size analysis**
 - ...and the questions we can formulate
- a **general method** to answer **sample size** questions
 - ...and obtain recommendation
- a software **implementation**
 - ...and an example
- end with a `/sample(=?size)/`



Simulation Approaches



- the process goes as follows:
 - select true parameter values for your model
 - generate one dataset with the true parameters
 - estimate the model parameters
 - test your hypothesis

- the process goes as follows:
 - select true parameter values for your model

- generate one dataset with the true parameters
- estimate the model parameters
- test your hypothesis

repeat many times

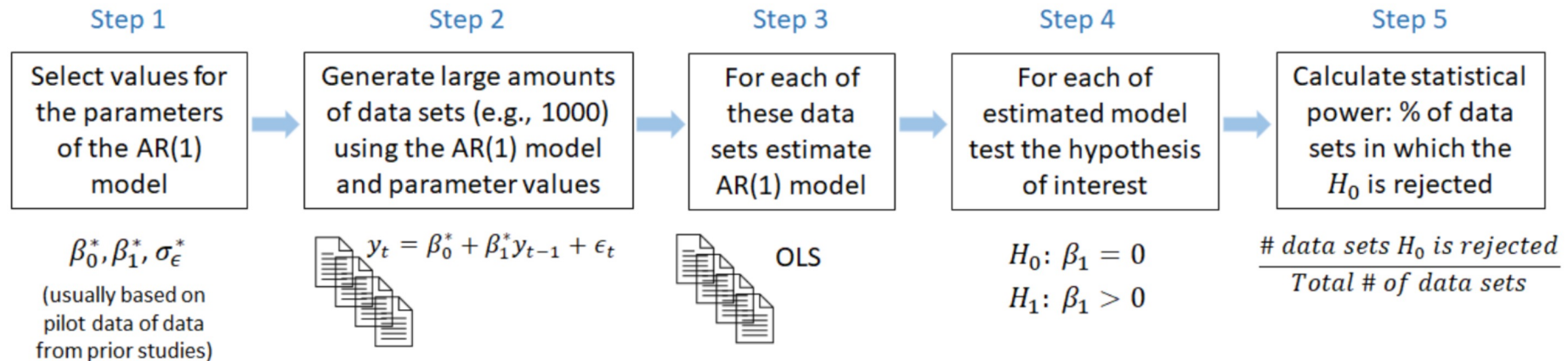
calculate
empirical power

We've seen this before...

Steps of the simulation-based approach

- Example: select the number of measurement occasions T to test if the autoregressive effect of PA is positive

Given T , Hypothesis of interest (e.g., $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 > 0$), and α



...and even how to code it

The Monte Carlo simulation function

This function conducts the Monte Carlo simulation for a set of sample sizes (i.e., several different number of observations) and computes the statistical power for a given hypothesis. It takes several arguments as follows:

- `vars` is the number of variables of the VAR(1) model
- `Tobs_list` is a list of numbers of repeated measurements (i.e., `Tobs`)
- `delta` the intercept matrix
- `psi` the transition matrix (which contains the auto-regressive and cross-regressive effects)
- `sigma` the variance-covariance matrix of the innovation
- `R` is the number of Monte Carlo replicates (e.g., 1000)
- `alpha` is the Type I error rate (or significance level of a test statistic)

```
# Function to conduct the Monte Carlo power simulation.
mc_power <- function(vars, Tobs_list, delta, psi, sigma, R, alpha) {
  # Prepare simulation storage.
  df_pow <- data.frame()

  # For each sample size in the list.
  for (i in 1:length(Tobs_list)) {
    # Extract the sample size.
    Tobs <- Tobs_list[i]

    # Print the progress.
    print(paste0("Power analysis for N = ", Tobs))

    # For each Monte Carlo replication.
    for (r in 1:R) {
      # Generate data.
      data <- sim_VAR_data(vars, Tobs, delta, psi, sigma)
```


We still leverage this this setup...

- but thinking about it more generally
- along two acts

We still leverage this this setup...

- but thinking about it more generally
- along two acts

the first act

What is the required input for running a simulation-based power analysis?

We still leverage this this setup...

- but thinking about it more generally
- along two acts

the second act

How can we process the input to get a sample size recommendations?

The Requirements

what



For a simulation approach we need to:

- generate or **specify true model** parameters
- **generate data** based on the true model parameters
- **estimate model** parameters from data
- specify a **performance measure** of interest
- specify a working **definition for power**



For a simulation approach we need to:

- generate or specify true model parameters
- **generate data** based on the true model parameters
- **estimate model** parameters from data
- specify a performance measure of interest
- specify a working definition for power



what to be able to perform



For a simulation approach we need to:

- generate or **specify true model** parameters
- generate data based on the true model parameters
- estimate model parameters from data
- specify a **performance measure** of interest
- specify a working **definition for power**

what to be able to provide



True Model Parameters

The Requirements

- the set of hypothesized model values used to generate data
- akin to an effect size in typical power analysis
- let's call it ω



- the **observed** data
 - is a sample from a data generating process with unknown parameters
- the **generated** data
 - is what we get when we pretend to know
 - the data generating process and
 - the values of its parameters → our hypothesized 🗨
 - we typically generate datasets of varying sizes for a given 🗨

- the observed data
 - is a sample from a data generating process with unknown parameters
- the **generated data**
 - is what we get when we provide
 - the data generating process
 - the values of its parameters → our hypothesized 🤖
 - we typically generate datasets of varying sizes for a given 🤖

What do you think we need the generated data for?



Estimated Model Parameters

The Requirements

- if Θ represents the hypothesized **true model** parameters
- then $\hat{\Theta}$ holds the **estimated model** parameters
 - estimated from the **generated data**

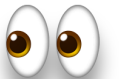


Estimated Model Parameters

The Requirements

- if Θ represents the hypothesized **true model** parameters
- then $\hat{\Theta}$ holds the **estimated model** parameters
 - estimated from the generated data

What does your intuition say will happen to $\hat{\Theta}$ if the generated dataset is very large?



- is a statement about the data generating process
 - quantifies the quality of the estimation
- expressed as $f(\Theta, \hat{\Theta})$ that
 - compares the **true model** parameters in Θ to the **estimated model** parameters in $\hat{\Theta}$
 - and the result of this comparison is dependent on the sample size

- is a statement about the data generating process
- expressed as $f(\Theta, \hat{\Theta})$ that
 - compares the true model parameters in Θ to the **estimated model** parameters in $\hat{\Theta}$
 - and the result of

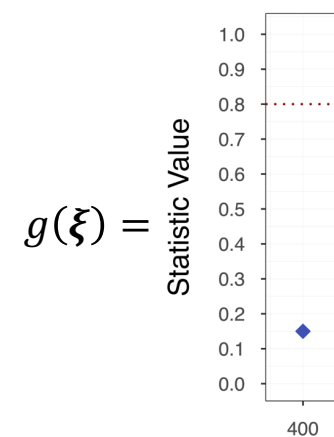
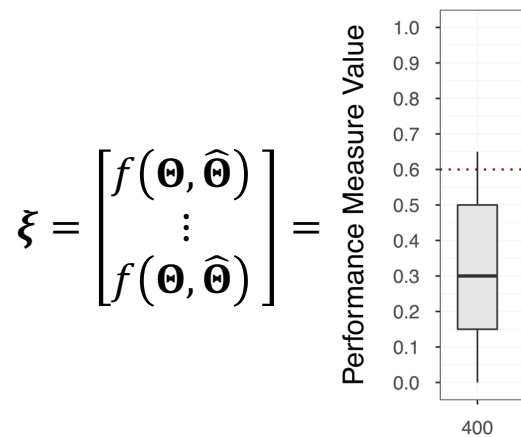
How is the performance measure $f(\Theta, \hat{\Theta})$
connected to the sample size?



- is a statement about the data generating process
- expressed as $f(\Theta, \hat{\Theta})$ that
 - compares the **true model** parameters in Θ to the **estimated model** parameters in $\hat{\Theta}$
 - and the result of this comparison is dependent on the sample size
- should be driven by the research question
- has a target value δ

- is a definition for the empirical power
 - that tells us how we want to observe the performance measure
 - e.g., we want a sample size such that **80%** of the performance measures reached the target δ

- is a definition for the empirical power
 - that tells us how we want to observe the performance measure
 - e.g., we want a sample size such that **80%** of the performance measures reached the target δ
 - is expressed as a function $g(\xi)$ with a target τ
 - where



The Required Input

The Requirements

what
in a nutshell



The Required Input

The Requirements

- true model 🗑️
 - it can be many things

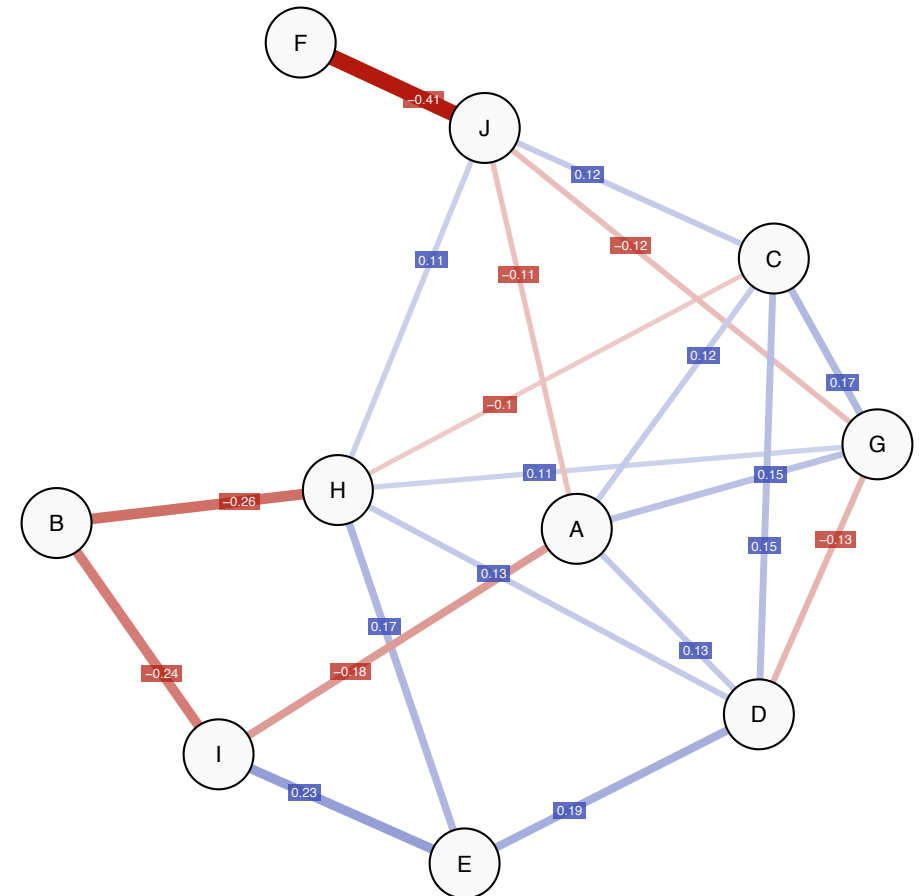


The Required Input

The Requirements

- true model 🗣️
 - it can be many things

Gaussian Graphical Model

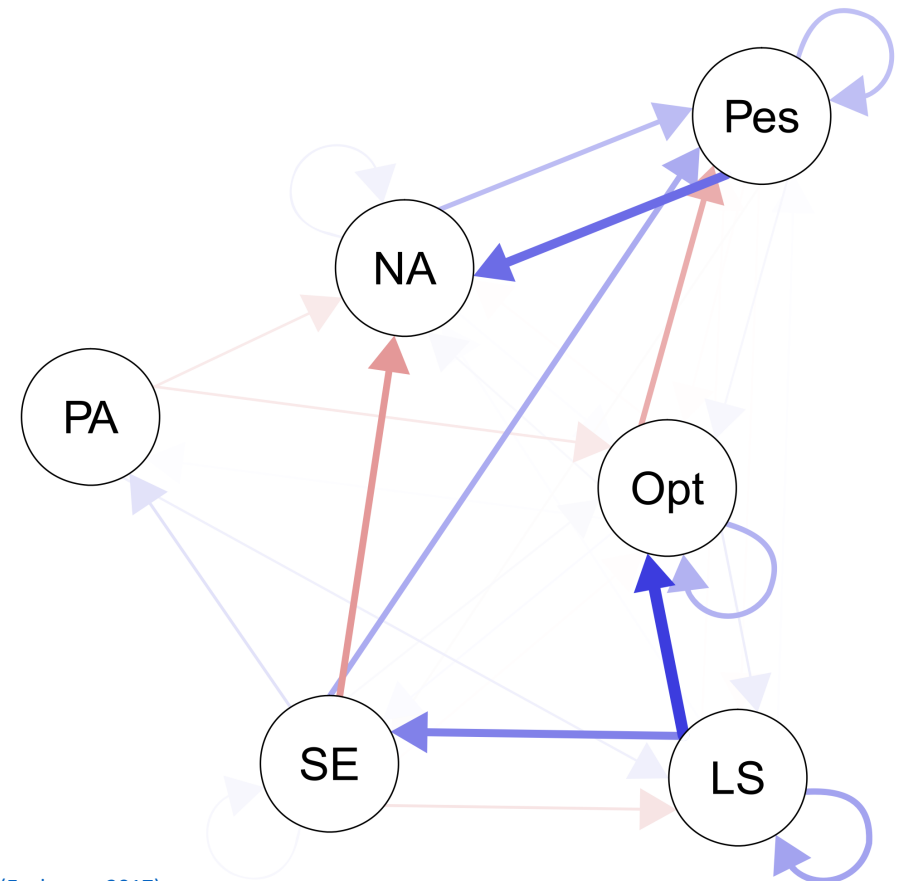


The Required Input

The Requirements

- true model 🗣️
 - it can be many things

Vector Autoregressive Model



(Epskamp, 2017)

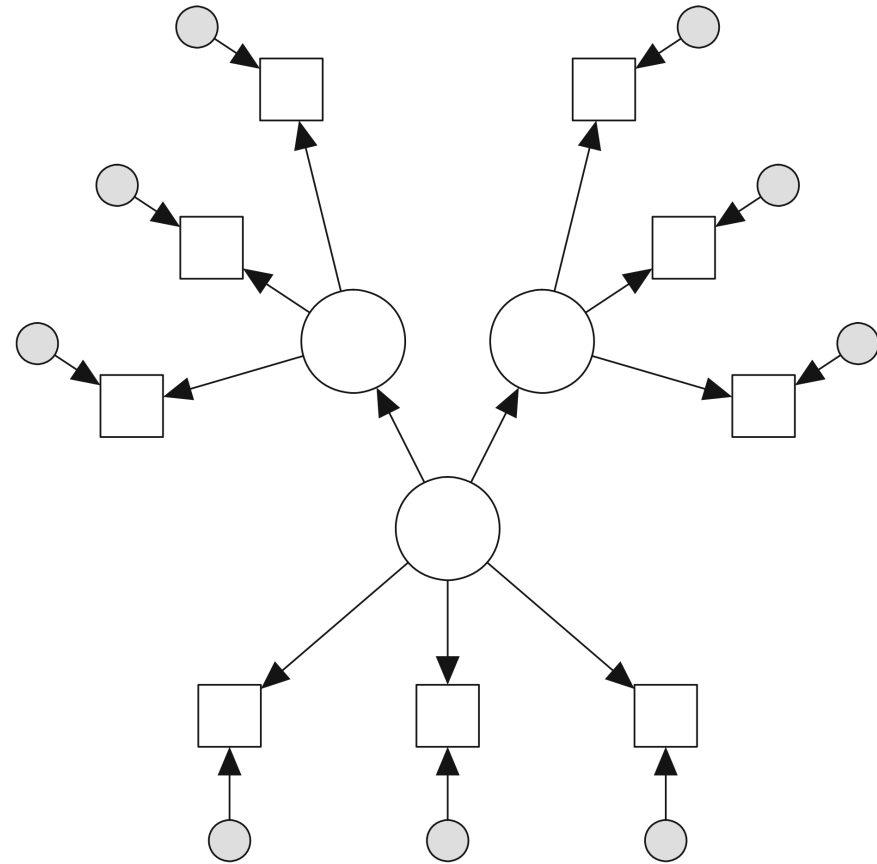


The Required Input

The Requirements

- true model 🗑️
 - it can be many things

Structural Equation Model



[\(Epskamp et al., 2017\)](#)



The Required Input

The Requirements

- true model Θ
- performance measure $f(\Theta, \hat{\Theta})$
 - should reflect the research question

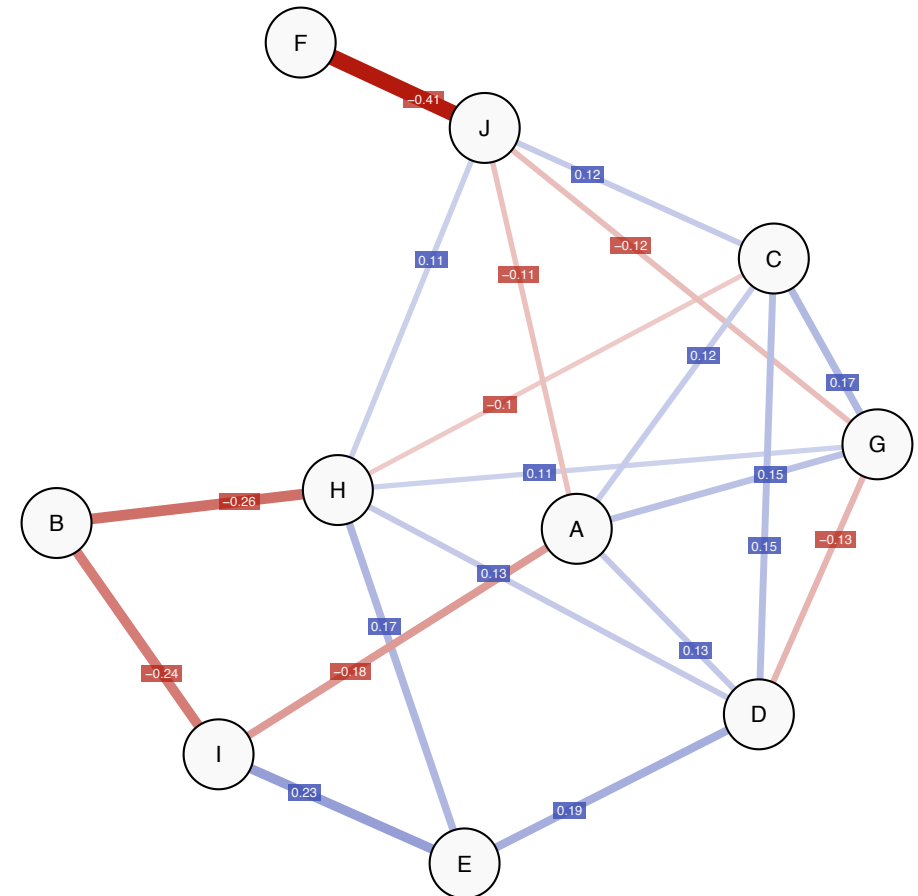


The Required Input

The Requirements

- true model Θ
- performance measure $f(\Theta, \hat{\Theta})$
 - should reflect the research question
 - e.g., suppose we want to recover the network structure
 - we look at **sensitivity** \rightarrow the proportion of edges correctly estimated to be present

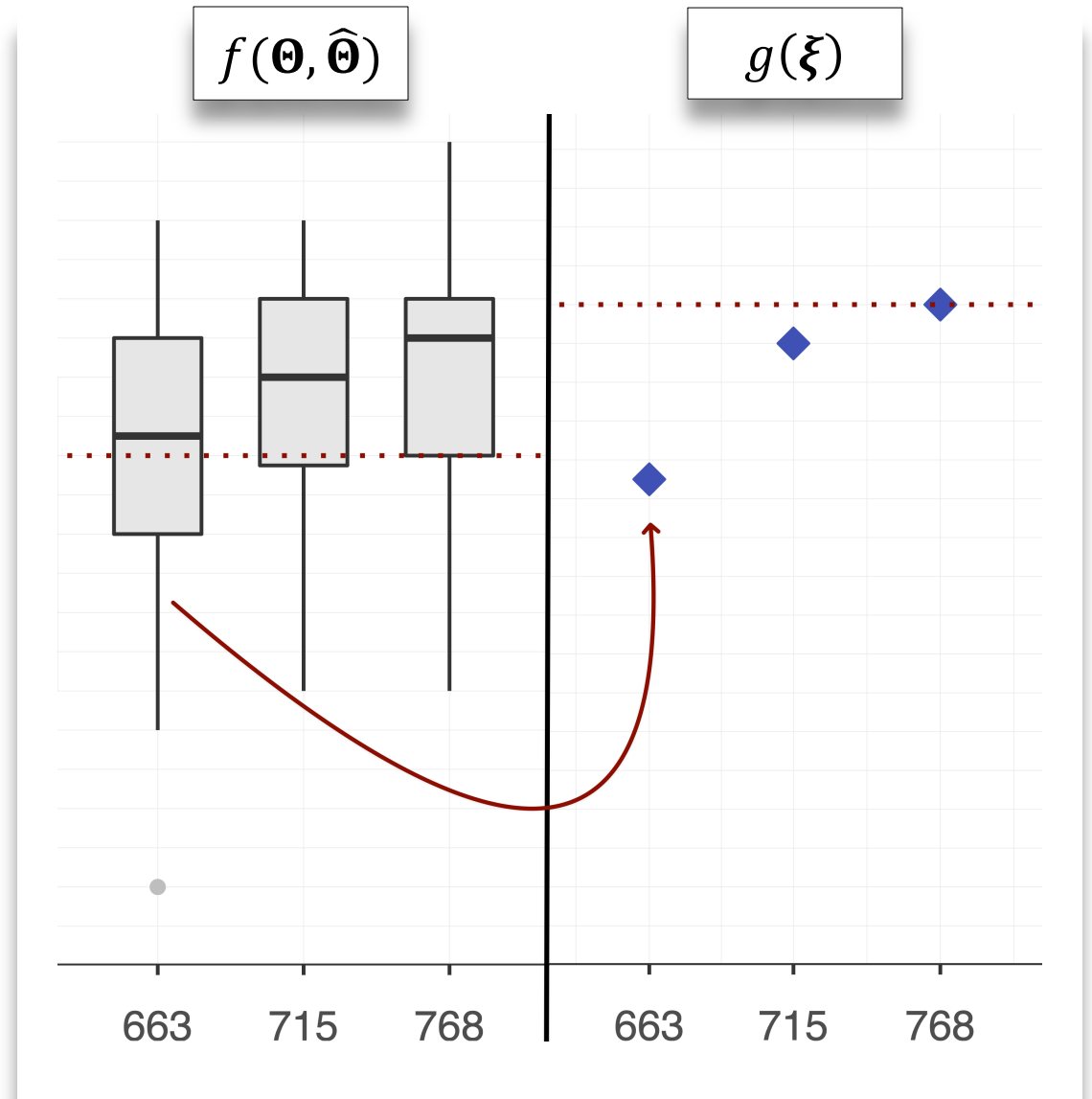
Gaussian Graphical Model



The Required Input

The Requirements

- true model Θ
- performance measure $f(\Theta, \hat{\Theta})$
- a statistic $g(\xi)$
 - most intuitively defined as a probability, but it may take other forms



The Required Input

The Requirements

- true model Θ
- performance measure $f(\Theta, \hat{\Theta})$
- a statistic $g(\xi)$



The Required Input

The Requirements

Based on this input, we can ask...

Given the hypothesized Θ , what sample size do we need to observe a $f(\Theta, \hat{\Theta}) \geq \delta$ with probability τ as defined by $g(\xi)$?



One could ask...

I have some idea about a $VAR(1)$ model I plan to fit, and I want to test that all my autoregressive coefficients are significant with a power of 0.8.
How much data do I need?

But in reality...

I have some idea about a $VAR(1)$ model I plan to fit, and I want to test that all my autoregressive coefficients are significant with a power of 0.8.

How much data do I need?

Here's The Deal

The Requirements



you provide us with the
required input



we provide you with the
sample size



The Method

how



We use a three-step Monte Carlo (MC) method that

- iteratively searches for an optimal sample size
- efficiently concentrates the MC simulations on relevant sample sizes
- can **extend** to other models and performance measures



The goal of this step is to get a rough understanding of the **behavior of**
 $f(\Theta, \hat{\Theta})$ as a **function of sample size**.

Step 1

The Method

- start with a candidate sample size range \mathbb{N}_S
- select T equidistant samples $\mathcal{S} = \{s_1, \dots, s_T\} \subseteq \mathbb{N}_S$

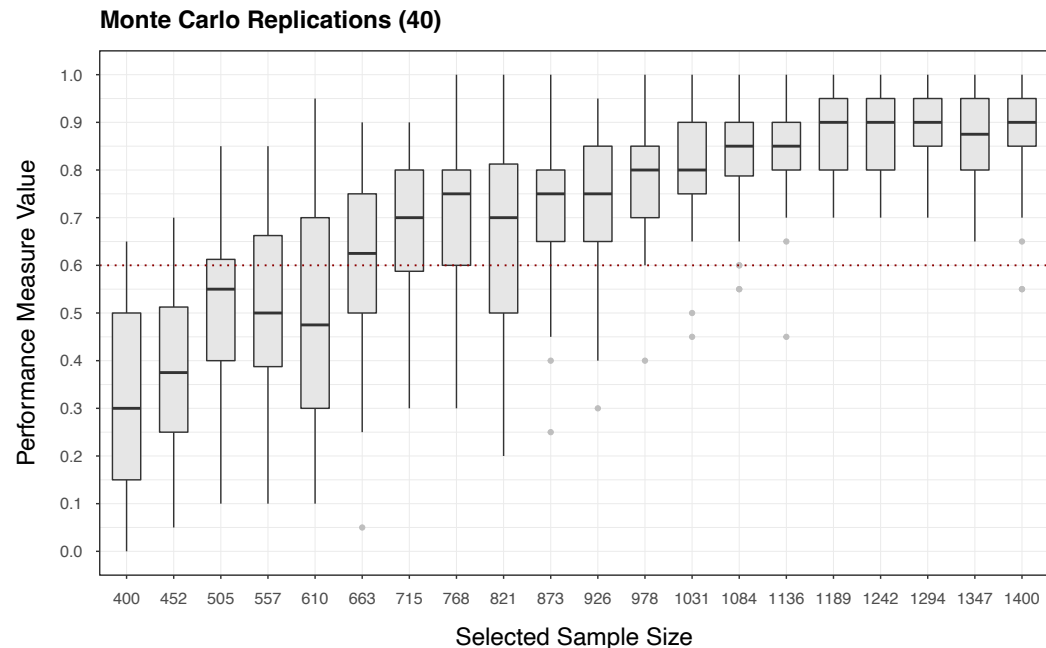


- start with a candidate sample size range \mathbb{N}_s
- select T equidistant samples $S = \{s_1, \dots, s_T\} \subseteq \mathbb{N}_s$
- for each $s_t \in S$ **perform R MC replications** as follows:
 - generate data with s_t number of cases using Θ
 - estimate $\hat{\Theta}$ using the generated data
 - compute $f(\Theta, \hat{\Theta})$

Step 1

The Method

- obtain $R \times T$ matrix \mathbf{E} , where each entry is a performance measure computed for a sample size during a MC replication

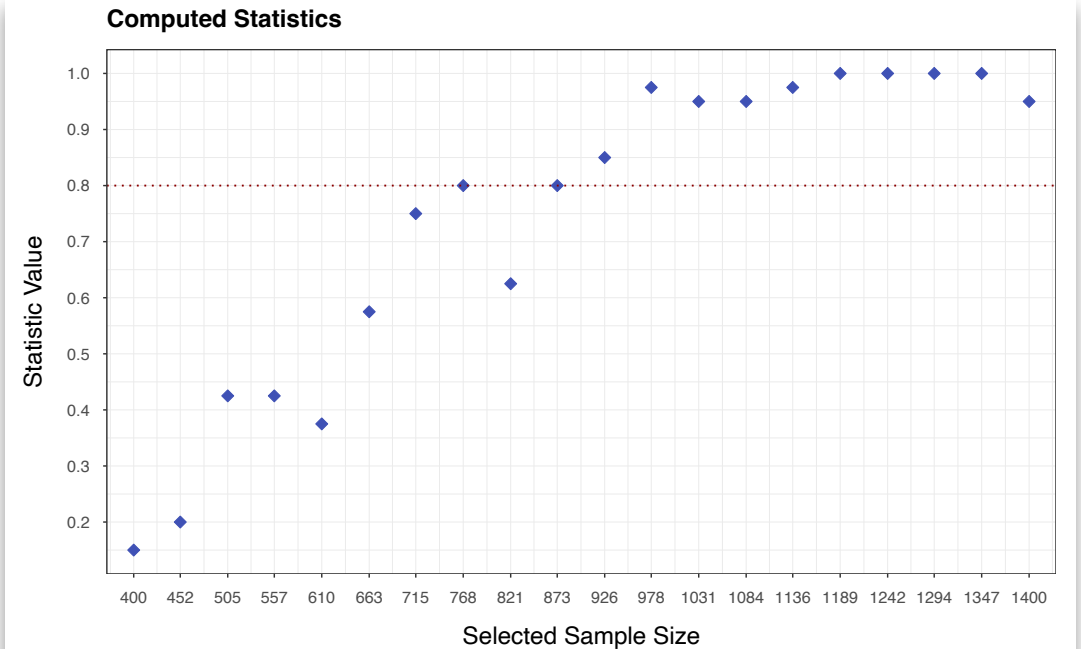
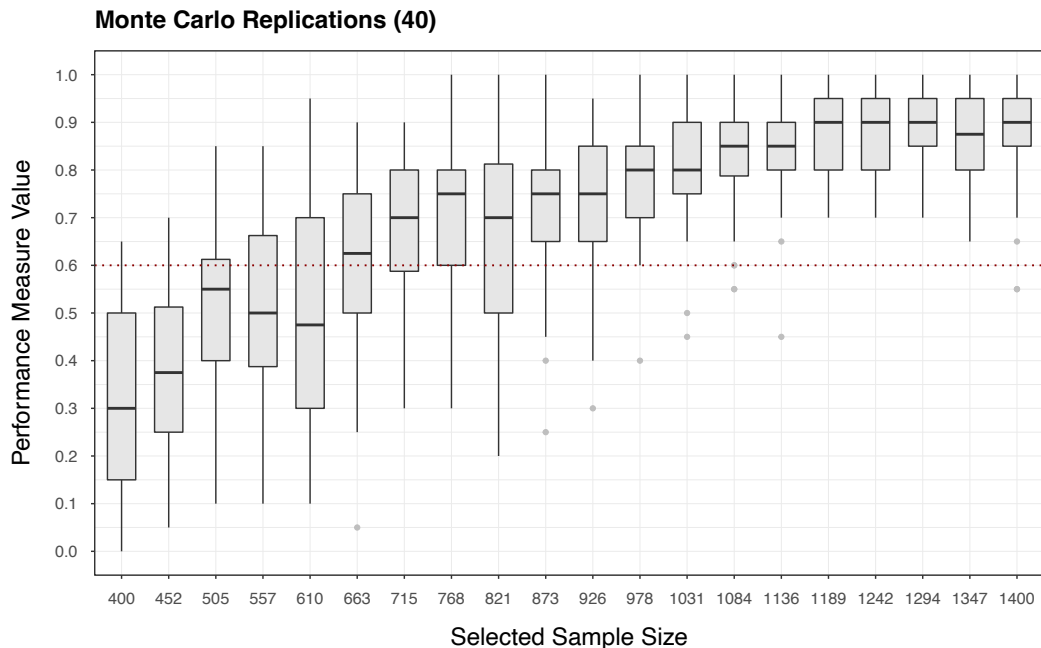


Step 1

The Method

- obtain $R \times T$ matrix Ξ , where each entry is a performance measure computed for a sample size during a MC replication

- apply $g(\xi)$ over each column of Ξ to **compute the statistic** (e.g., power)

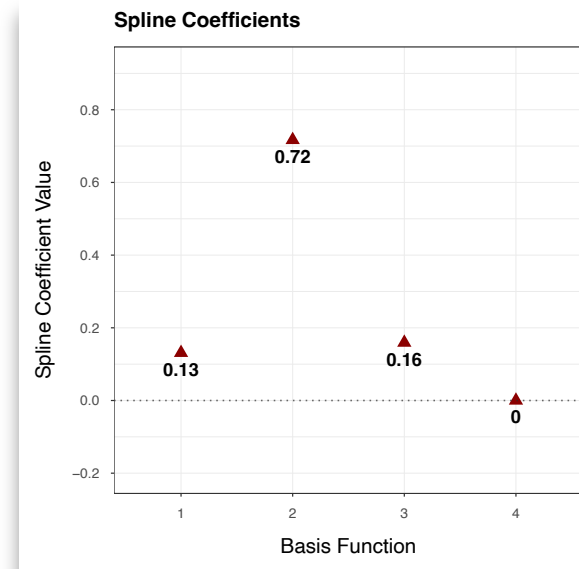
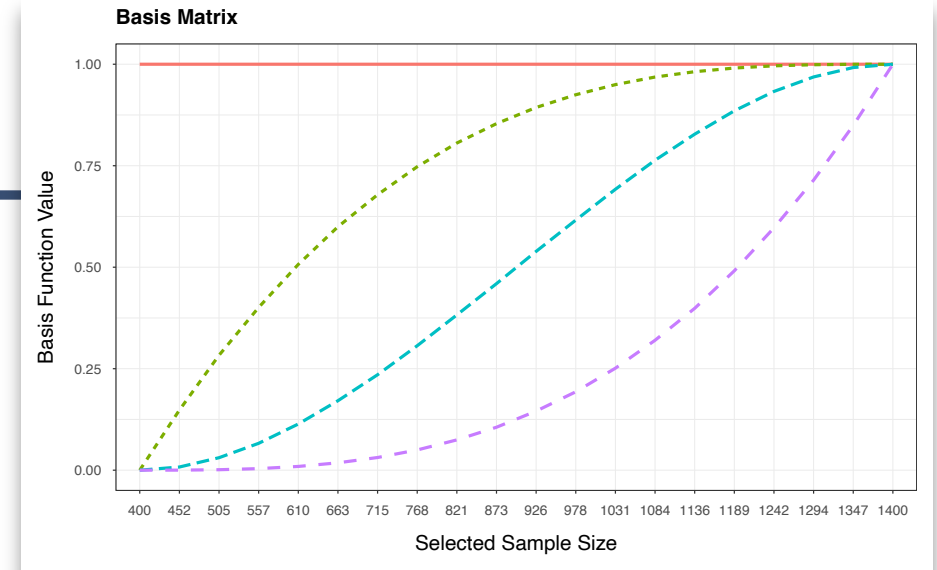


The goal of this step is to obtain a smooth (power) function and **interpolate the statistic** for **all sample sizes** in the range \mathbb{N}_s .

Step 2

The Method

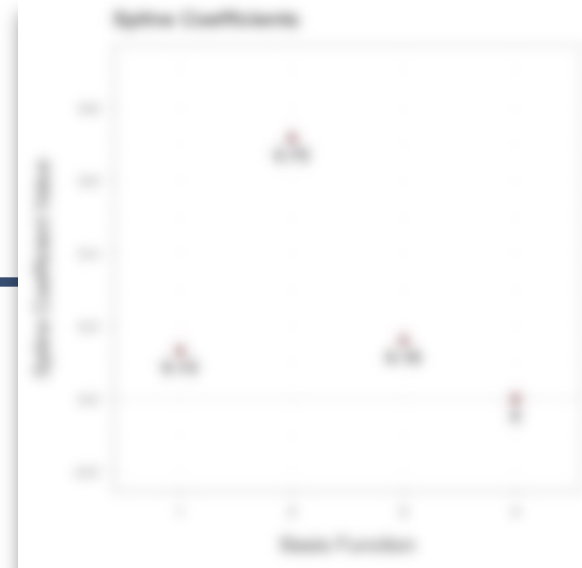
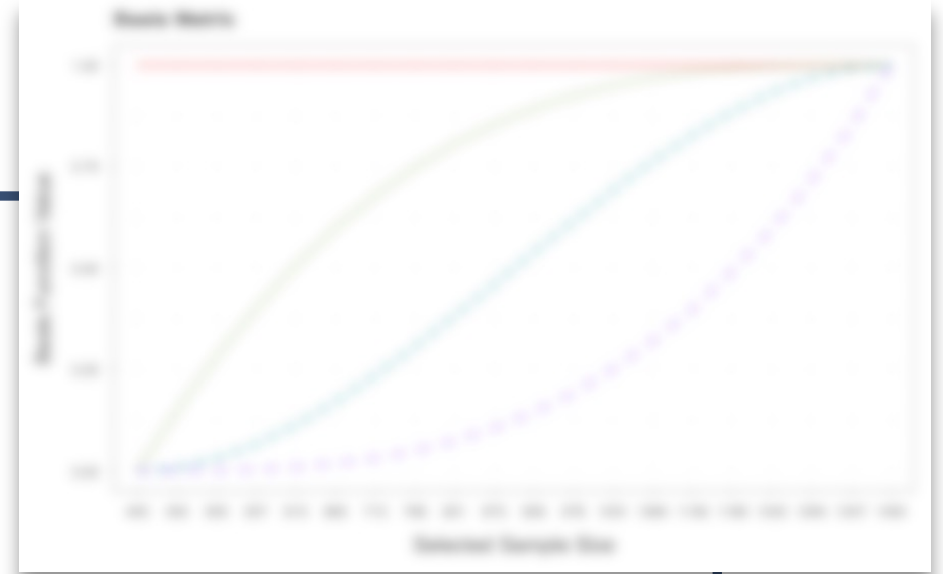
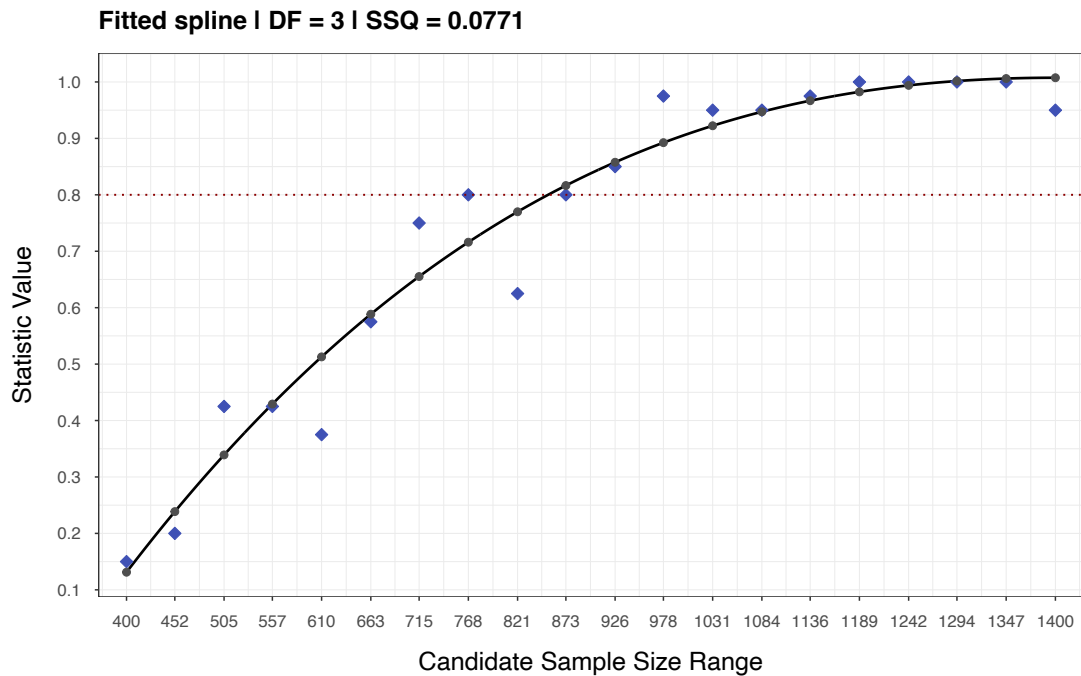
- **assume monotonicity** and use cubic *I-spline* bases with inner knots selected based on cross-validation



Step 2

The Method

- **assume monotonicity** and use cubic *I*-Spline bases with inner knots selected based on cross-validation



The goal of this step is to **account for the MC error** and provide a measure of uncertainty around the interpolated spline.

Step 3

The Method

- use **stratified bootstrapping** to represent the variability in the replicated performance measures for each sample size $s_t \in S$

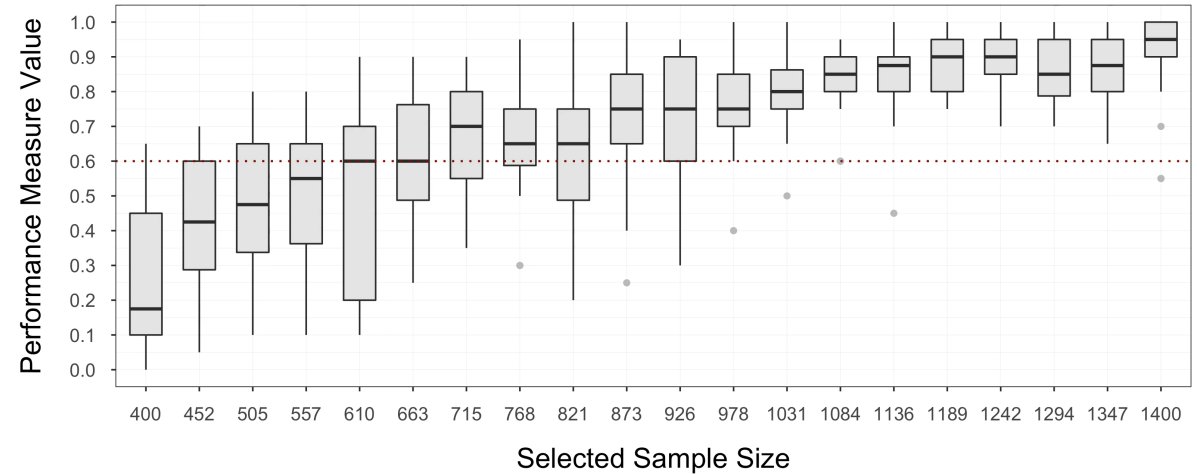


Step 3

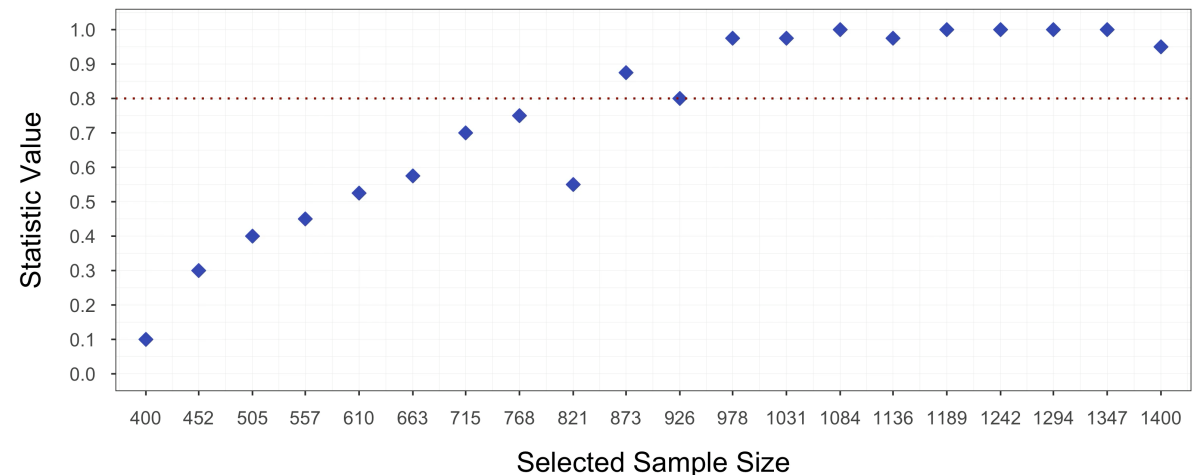
The Method

- use stratified bootstrapping to represent the variability in the replicated performance measures for each sample size $s_t \in S$
- we bootstrap the **performance measures** and, thus, re-estimating the model is not necessary

Monte Carlo Replications (40)



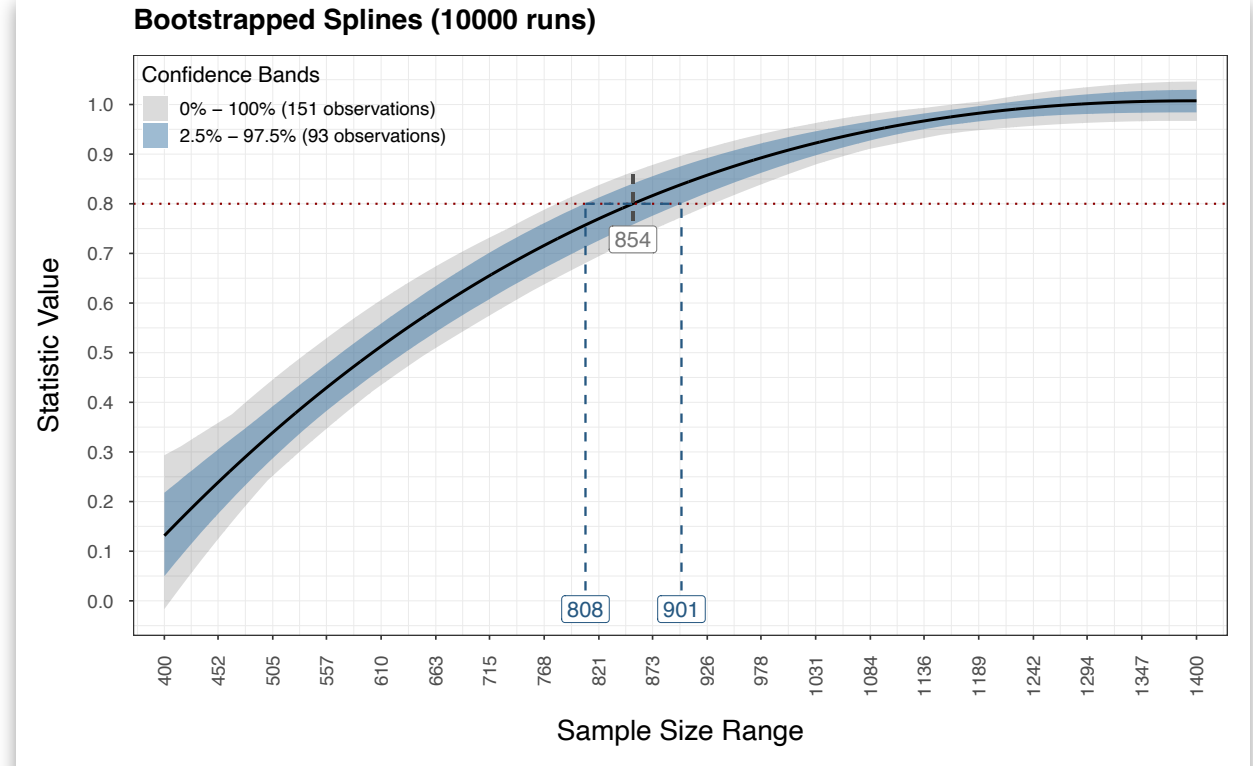
Computed Statistics



Step 3

The Method

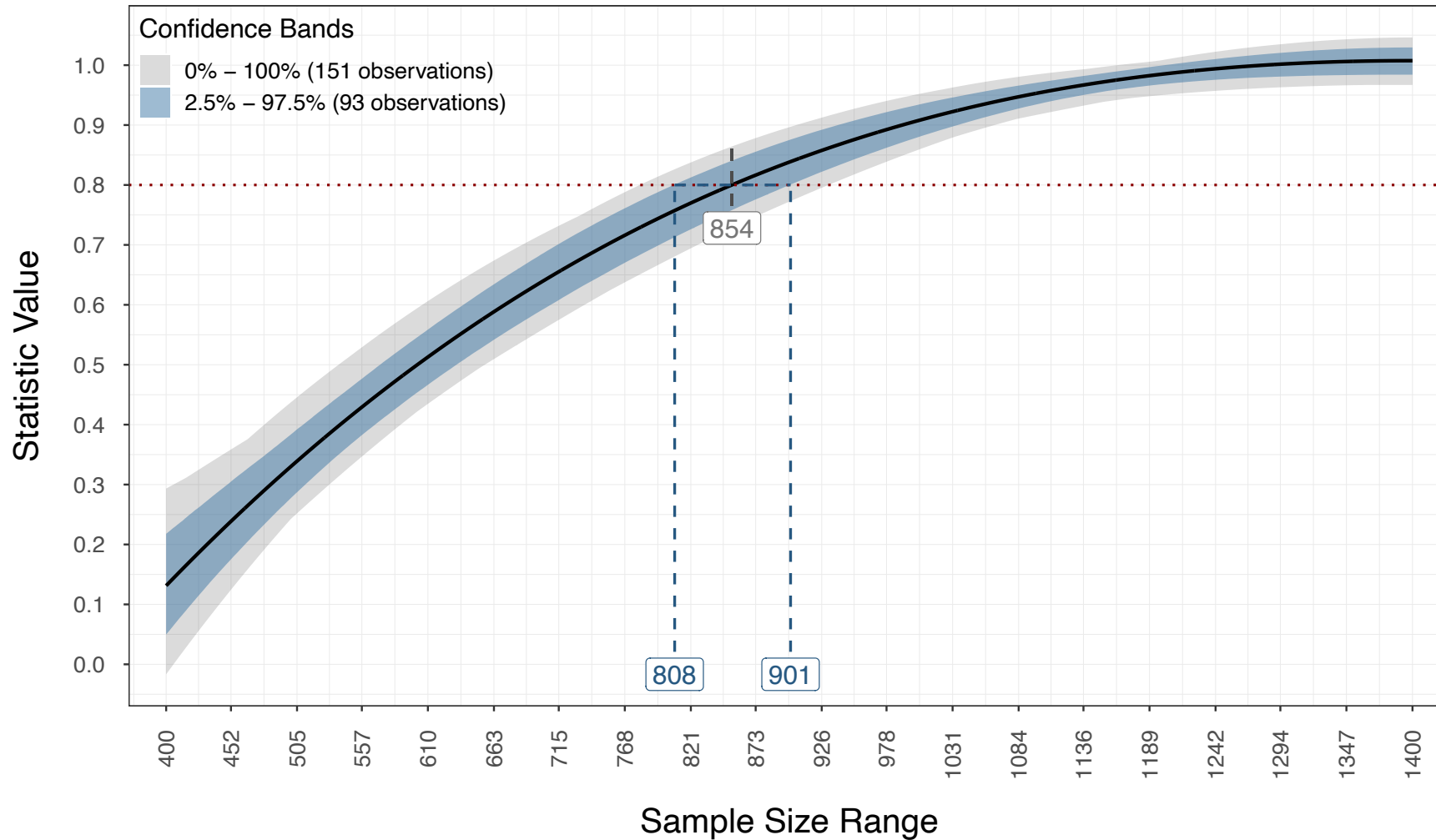
- use stratified bootstrapping to represent the variability in the replicated performance measures for each sample size $s_t \in S$
- we bootstrap the performance measures and, thus, re-estimating the model is not necessary
- fit a new spline to each bootstrapped matrix of performance measures



Step 3

The Method

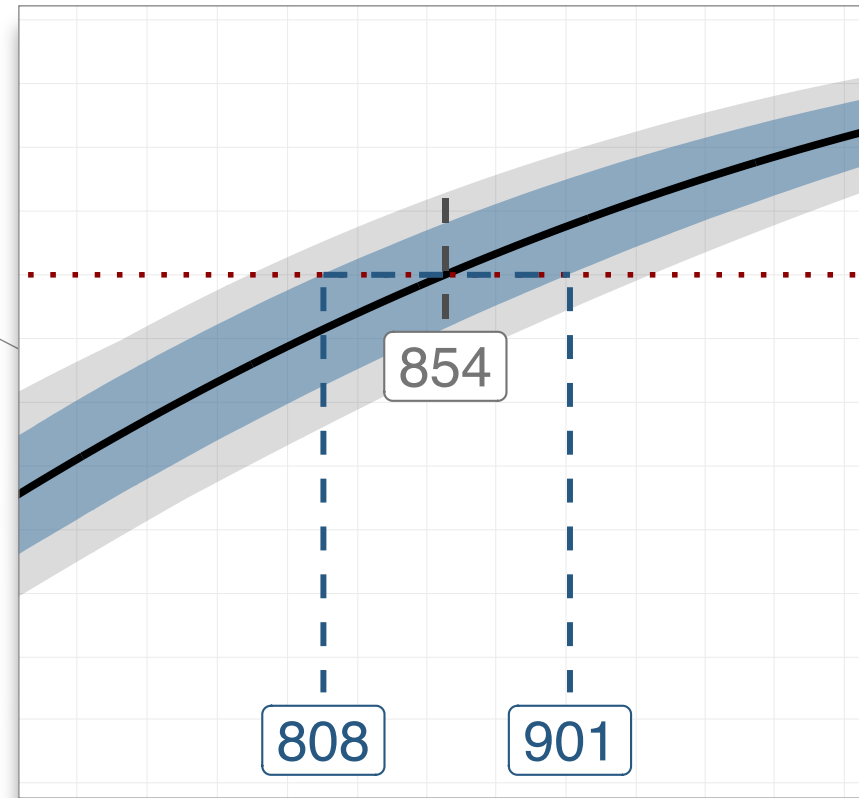
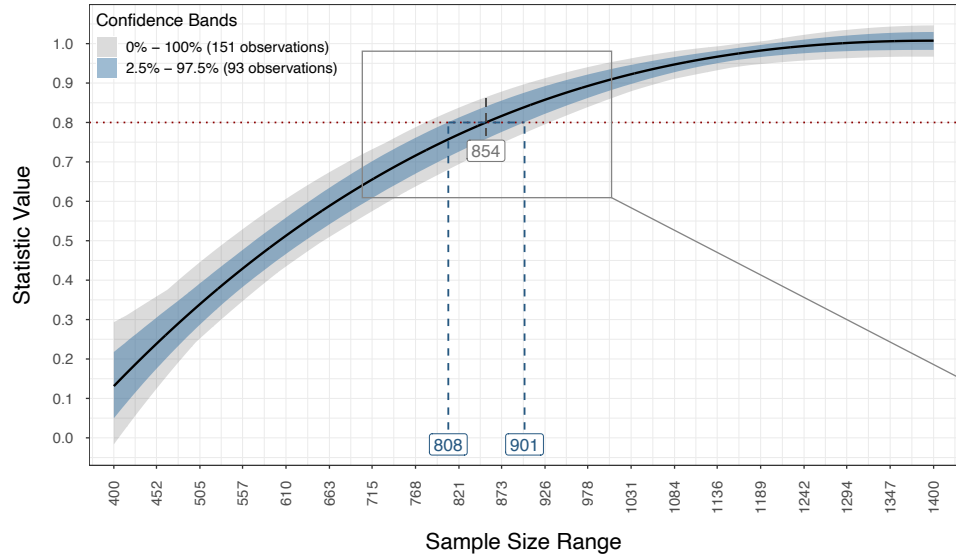
Bootstrapped Splines (10000 runs)



Step 3

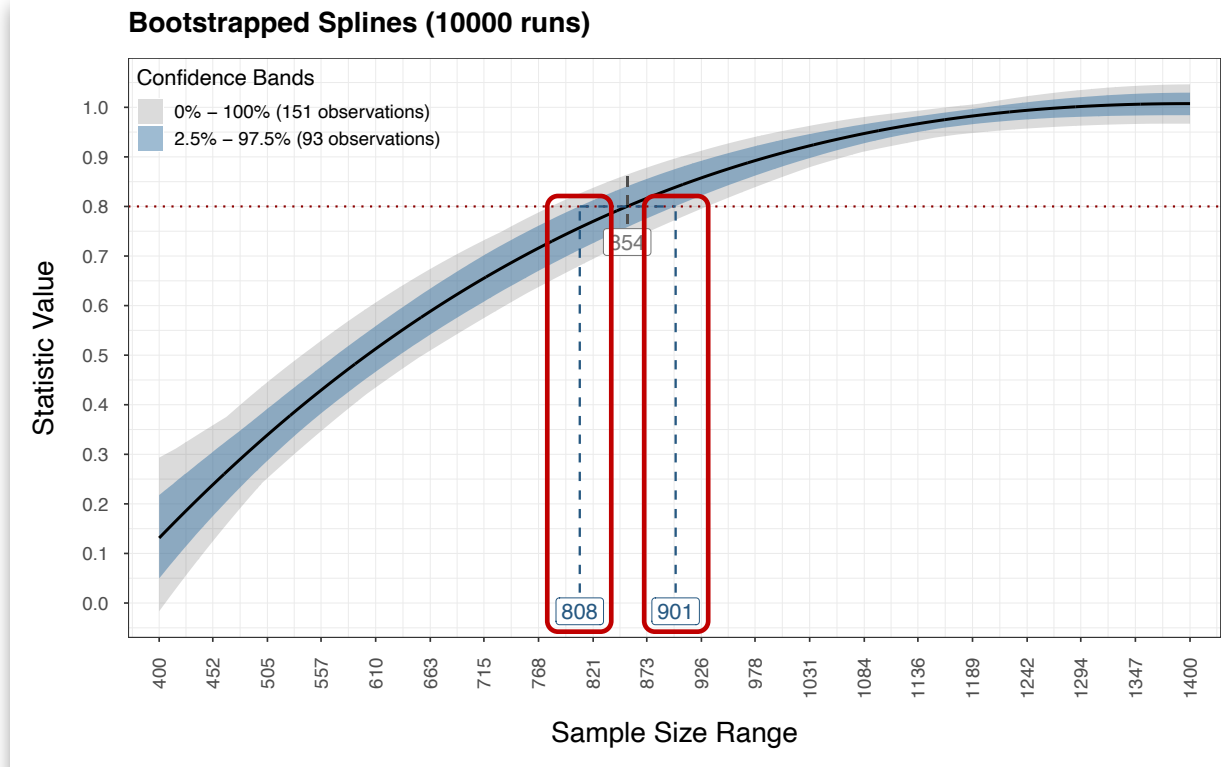
The Method

Bootstrapped Splines (10000 runs)



Convergence

The Method



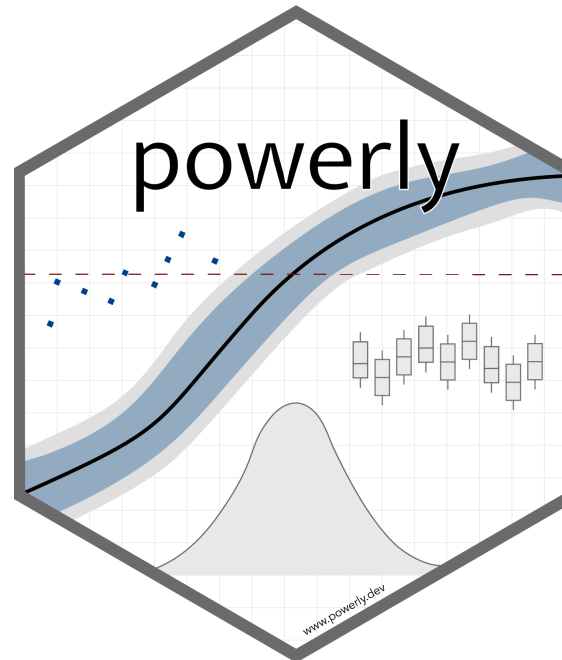
- **update candidate range** N_s based on the confidence bands
- **repeat** Steps 1 to 3 until range N_s becomes small enough



The Implementation



- an **R** package



powerly.dev

In Code

The Implementation

```
# Load the library.  
library(powerly)
```



```
# Load the library.  
library(powerly)
```

```
# Generate a true model.  
true_model <- generate_model(  
  type = "...",  
  ...  
)
```

```
# Load the library.  
library(powerly)
```

```
# Generate a true model.  
true_model <- generate_model(  
  type = "...",  
  ...  
)
```

```
# Load the library.  
library(powerly)
```

```
# Generate a true model.  
true_model <- generate_model(  
  type = "...",  
  ...  
)
```

```
# Run the method.  
results <- powerly(  
  range_lower = 300,  
  range_upper = 1000,  
  samples = 30,  
  replications = 20,  
  measure = "...",  
  statistic = "power",  
  measure_value = .6,  
  statistic_value = .8,  
  model = "...",  
  model_matrix = true_model  
)
```

```
# Load the library.  
library(powerly)
```

```
# Generate a true model.  
true_model <- generate_model(  
  type = "...",  
  ...  
)
```

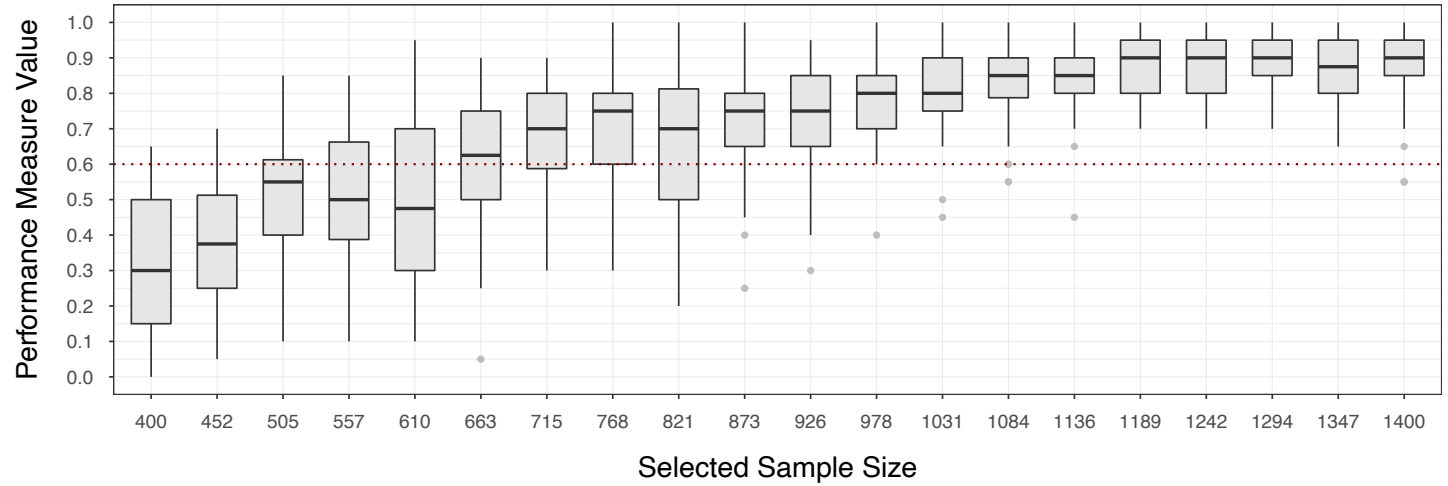
```
# Run the method.  
results <- powerly(  
  range_lower = 300,  
  range_upper = 1000,  
  samples = 30,  
  replications = 20,  
  measure = "...",  
  statistic = "power",  
  measure_value = .6,  
  statistic_value = .8,  
  model = "...",  
  model_matrix = true_model  
)
```

Step 1

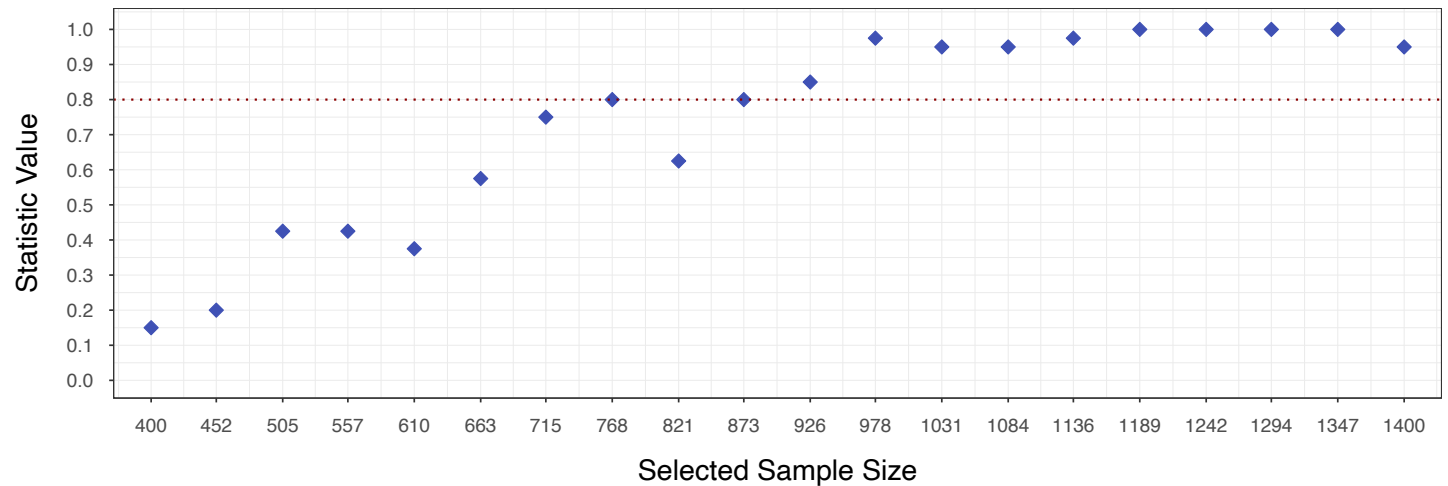
The Implementation

```
plot(results, step = 1)
```

Monte Carlo Replications (40)



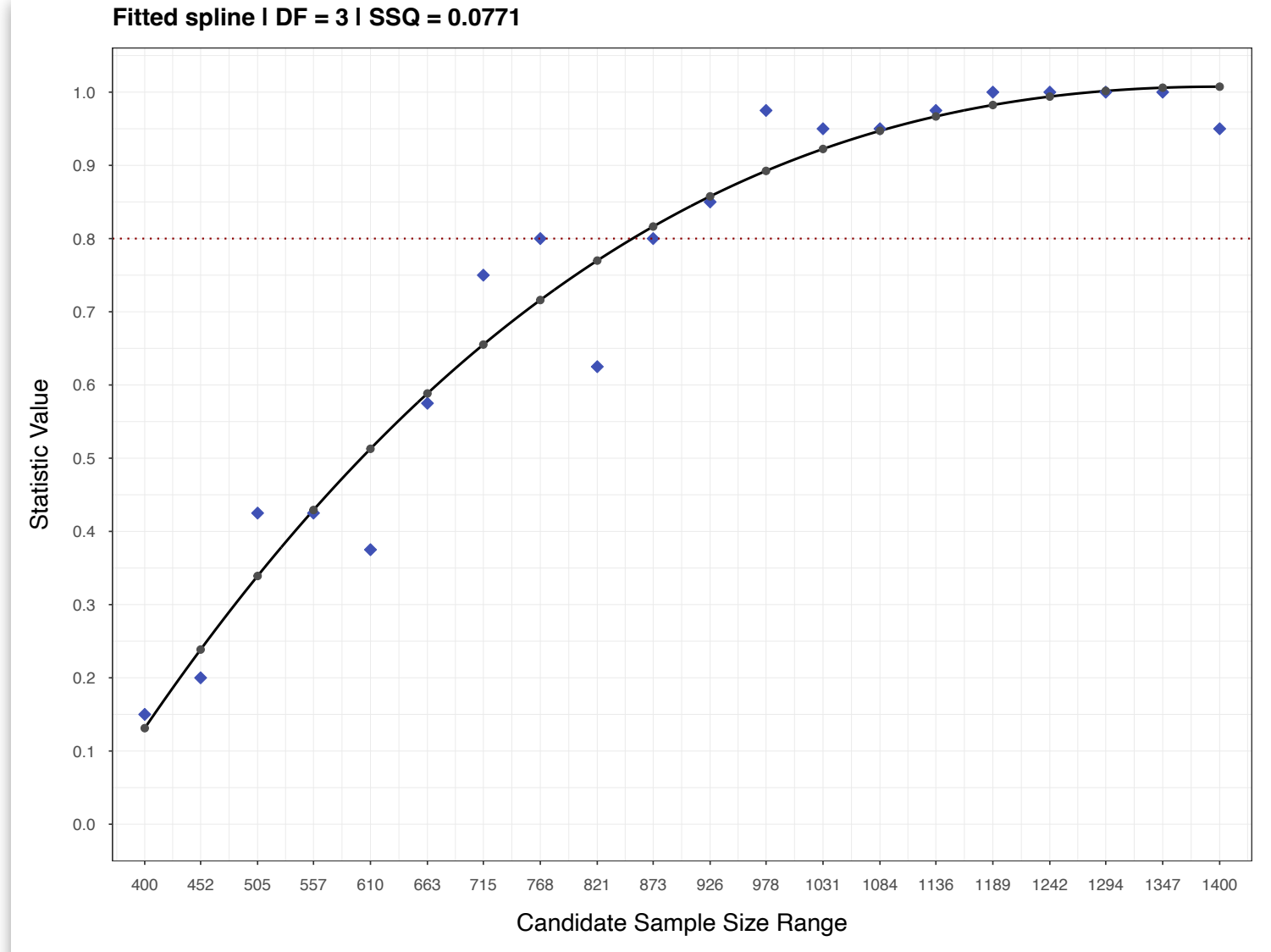
Computed Statistics



Step 2

The Implementation

```
plot(results, step = 2)
```

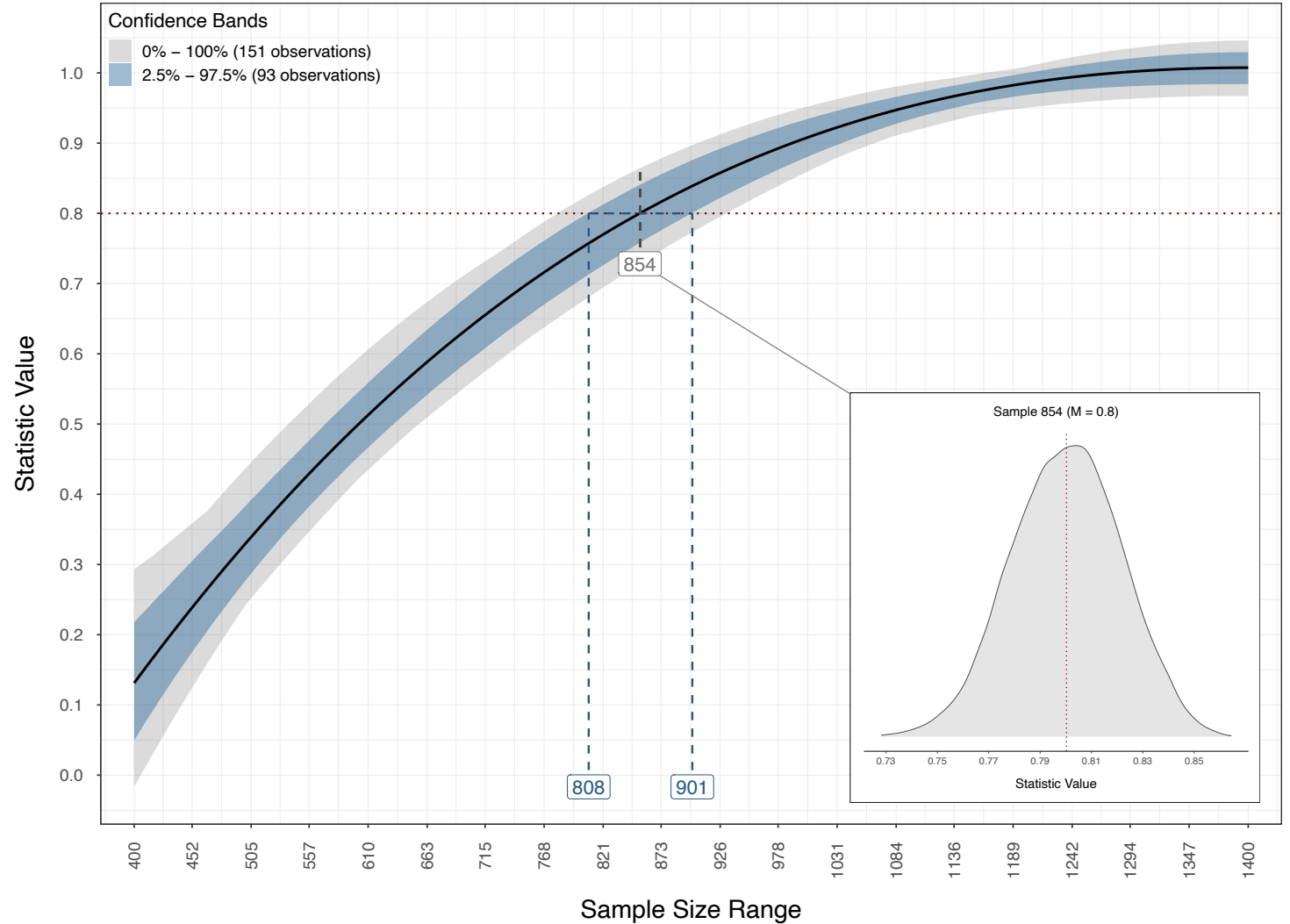


Step 3

The Implementation

```
plot(results, step = 3)
```

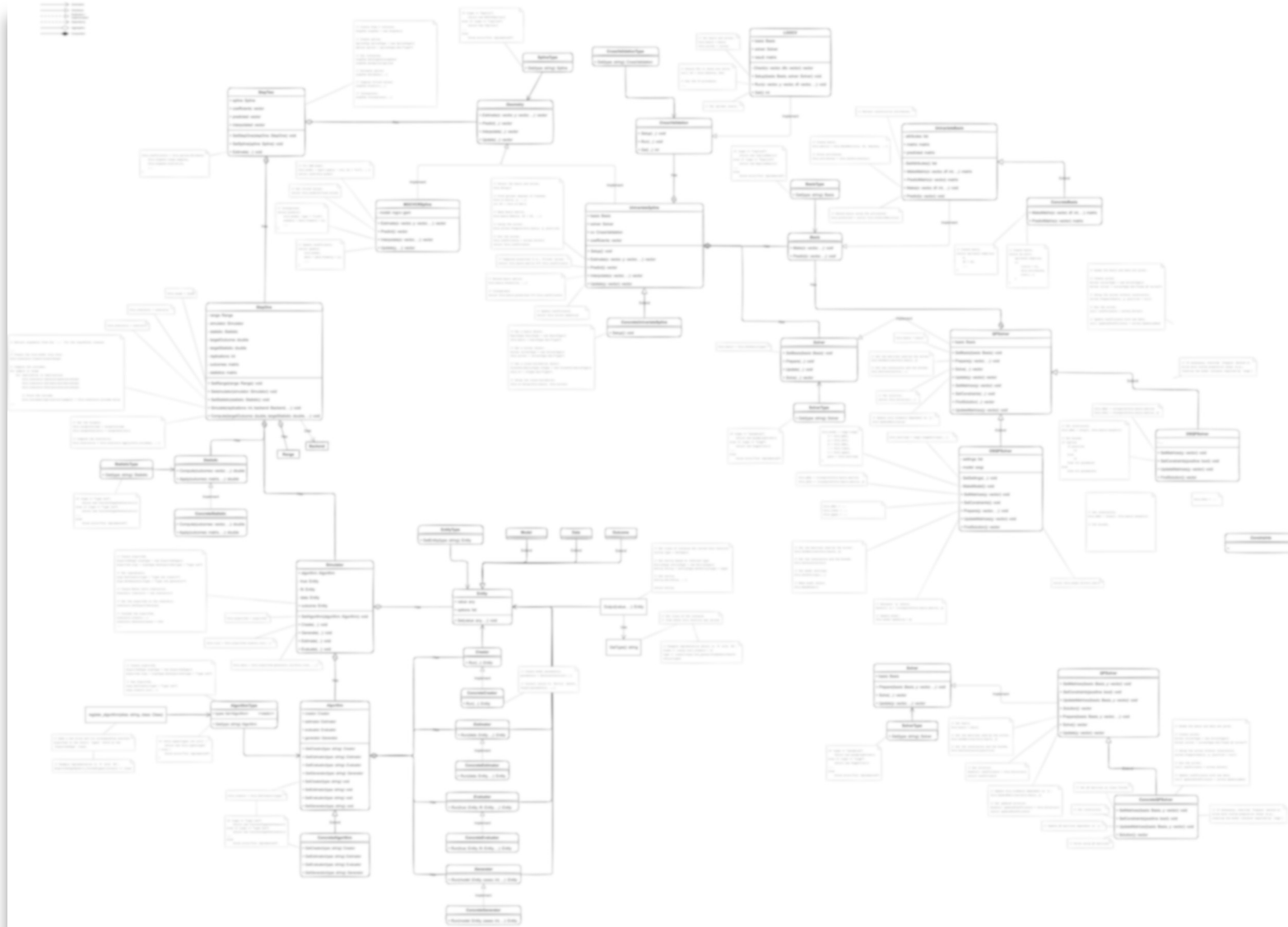
Bootstrapped Splines (10000 runs)



Positive Lookahead



A General Framework



by [Marie Mainguy](#)



Why?

- sample sizes tailored to specific research questions
- sample size analysis as an **ecosystem**
 - growing **collection** of models and performance measures
 - developer **API** for enabling sample size computations
- upcoming tutorial paper where we
 - discuss these ideas
 - and show how to apply them



Our Final Frontier

we aim to make sample size analysis so accessible that
there is no way around not doing it



by [Marie Mainguy](#)



samplesize.help