

# Appendix: Validation of power results

June 27, 2021

## Contents

<b>TODO</b>	<b>1</b>
<b>Introduction</b>	<b>1</b>
Monte Carlo Simulations . . . . .	1
Validation scenarios . . . . .	2
Simulation parameters . . . . .	2
<b>Validation results</b>	<b>4</b>
<b>Westfall-Young procedures</b>	<b>6</b>

## TODO

- Potentially problematic graphs
  - blocked\_c2\_3r  $R^2 = 0.6$
  - blocked\_c2\_3r  $ICC = 0.7$
  - blocked\_i1\_3r  $ICC = 0.7$
- Write up westfall-young
- Run all MDES and SS calculations
- Try additional simulation variations

## Introduction

This appendix discusses our work to validate that the power estimation methods work as intended.

We compare three different methods of estimating power:

- PUMP
- PowerUpR (only comparable for D1individual, unadjusted power)
- Monte Carlo Simulations

We compare point estimates of power from PUMP and PowerUpR for D1 individual, unadjusted power estimates. For all other types of power definitions and adjustments, we are only able to compare PUMP to the estimated power from Monte Carlo simulations. The simulations produce a 95% confidence interval for power, and we check that the PUMP estimate is within the confidence interval.

## Monte Carlo Simulations

The main work of this validation step was to design monte carlo simulations in order to estimate power. In order to estimate power by simulation, we follow the following steps.

For iteration  $s = 1, \dots, S$ :

1. Generate simulated data according to the assumed data generating process (DGP)
2. Generate simulated treatment assignment

3. Calculate p-value given simulated data, treatment assignment, and assumed model

At the end, a 95% confidence interval for power is calculated, assuming a conservative standard error estimate of  $\sqrt{0.25/\bar{S}}$ .

We also validate MDES and sample size calculations. For MDES, we choose one default scenario for each design and model, then input the already-calculated D1 individual power and see if the output MDES is the same as the original input MDES. Similarly, for sample size validation, we input the already-calculated D1 individual power and see if the output sample size (either  $J$  or  $K$  depending on design) is the same as the original sample size.

## Validation scenarios

We use the following adjustment procedures:

- Bonferroni
- Benjamini Hochberg (BH)
- Holm
- Westfall-Young Single Step (WY-SS)
- Westfall-Young Step Down (WY-SD)

We calculate power under the following definitions:

- Individual power for each outcome ( $M = 3$ ): D1indiv, D2indiv, D3indiv
- Mean individual power
- Minimum power: min1, min2
- Complete power

We consider the following designs:

- Blocked individual randomization, 2 level
  - constant effects (blocked\_i1\_2c)
  - fixed effects (blocked\_i1\_2f)
  - random effects (blocked\_i1\_2r)
- Blocked individual randomization, 3 level
  - random effects (blocked\_i1\_3r)
- Cluster randomization, 2 level
  - random effects (simple\_c2\_2r)
- Cluster randomization, 3 level
  - random effects (simple\_c3\_3r)
- Blocked cluster randomization, 3 level
  - fixed effects (blocked\_c2\_3f)
  - random effects (blocked\_c2\_3r)

## Simulation parameters

In order to validate that the method works in a wide range of scenarios, which vary the following parameters.

Parameters that vary:

Parameter	Default	Comparison values
school size $\bar{n}$	50	75, 100
$R^2$	0	0.6
$\rho$	0.5	0.2, 0.8
ATE (ES) true positives	(0.125, 0.125, 0.125)	(0.125, 0. 0)
ICC	0.2	0.7
$\omega$	0.1	0.8

We do not vary:

- $M = 3$
- $J$  and  $K$  are fixed for each scenario
- Scalar grand mean  $\Xi_0$ 
  - Correlations between school random effects and impacts  $\kappa$
  - $\rho$  informs all correlations; we keep the same correlation between covariates, residuals, impacts, random effects for all levels and across all outcomes

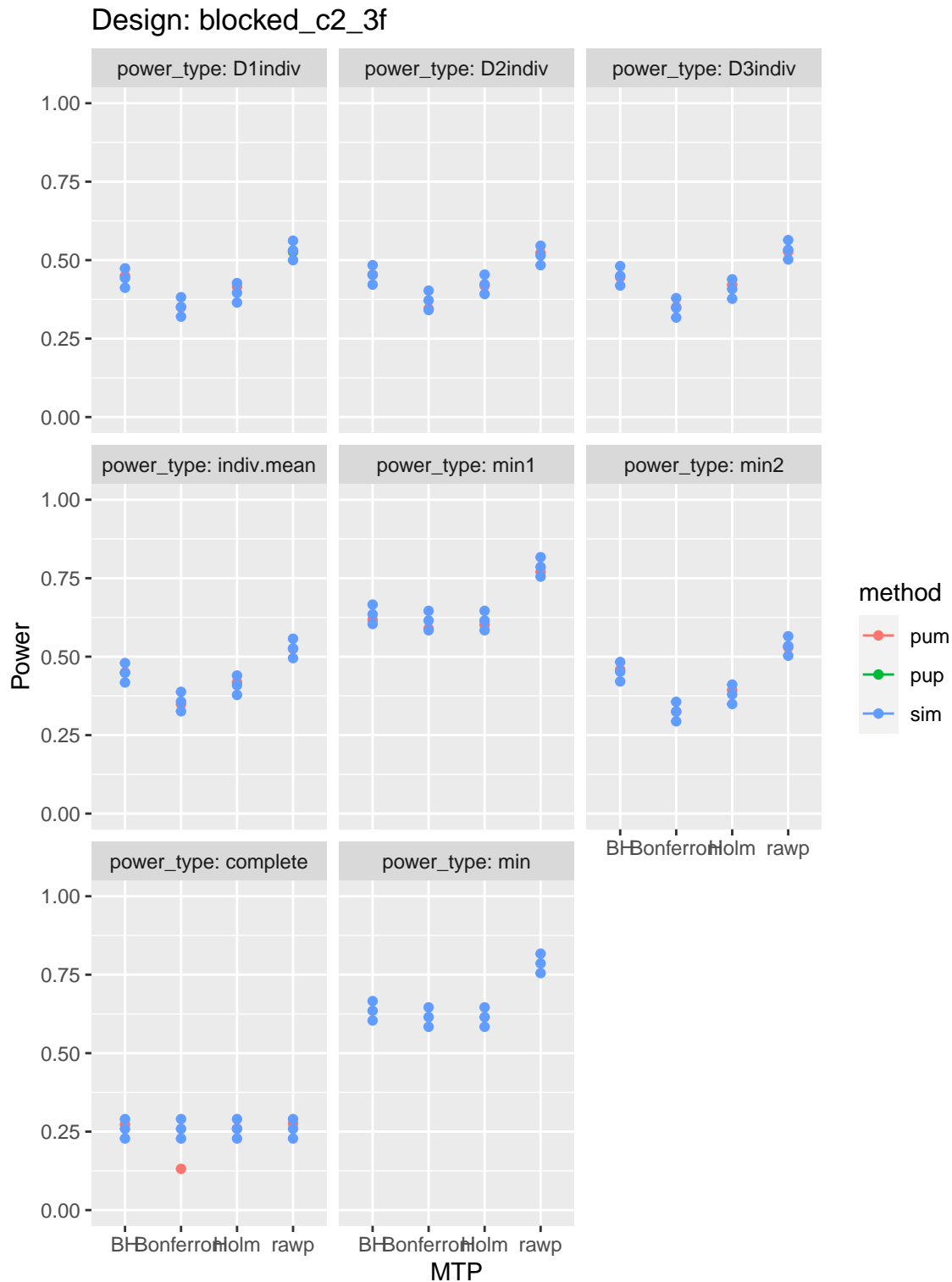
TODO: possible things we may want to vary:

- Correlations between school random effects and impacts  $\kappa$
- Combinations of parameters instead of one at a time
- A realistic set of parameters
- Extreme values, i.e.  $\rho = 0.95$

## Validation results

Below is an example of a graph we use for validation. The red dot is the PUM estimate of power, the green dot is the PowerUpR estimate of power, and the 95% confidence intervals based on the monte carlo simulations are shown in blue. To validate that PUMP produces the expected result, we want to see the red and green points match, and for the red point to be within the blue intervals. The plot shows the results across different types of power and different MTPs.

This graph in particular is for a blocked cluster design with fixed effects with a particular set of parameters. We repeat this graph for each set of parameters for each design.



Next, we validate MDES and sample size calculations. The first column shows the calculated MDES or sample size, the middle column is the power we plugged into the calculation, and the last column shows the MDES or sample size that we are targeting. Thus, ideally we want the first and last columns to match.

```
##
##
## +-----+-----+-----+-----+-----+-----+
```

```

## |      MTP      | Adjusted MDES | Diindiv Power | Target MDES |
## +=====+=====+=====+=====+
## |      rawp      |      0.1262      |      0.5202      |      0.125      |
## +-----+-----+-----+-----+
## | Bonferroni |      0.125      |      0.3405      |      0.125      |
## +-----+-----+-----+-----+
## |      BH      |      0.1247      |      0.4308      |      0.125      |
## +-----+-----+-----+-----+
## |      Holm      |      0.1246      |      0.3924      |      0.125      |
## +-----+-----+-----+-----+
##
## Table: blocked_c2_3f
##
##
## +-----+-----+-----+-----+-----+
## | MTP | Sample type | Sample size | Diindiv power | Target sample size |
## +=====+=====+=====+=====+=====+
## | BH | K | 11 | 0.438 | 10 |
## +-----+-----+-----+-----+-----+
## | Holm | K | 11 | 0.3977 | 10 |
## +-----+-----+-----+-----+-----+
##
## Table: blocked_c2_3f

```

## Westfall-Young procedures

The Westfall-Young procedure was validated separately due to its unique complications and computational burden.

First, we wrote a series of careful unit tests ensuring that our code worked as expected for each step of the WY procedure. One of these tests compared results from our procedure to the WY procedure in the multtest package, and found it matched quite well.

Second, we also tested the WY procedures using the same simulation procedure described above. For constant effect and fixed effect models, the power estimation matches between PUM and the simulations. However, for random effects models, the two methods diverge in some cases. We find that the simulated power matches the PUMP-calculated power only when (1) there is a large number of blocks/clusters, and (2) the user has a large number of WY permutations. Below, we discuss the single-step procedure because it is simpler, although the same concepts hold for the step-down procedure. Although it is difficult to verify why this discrepancy occurs, our hypothesis is that this behavior occurs due to the combination of the sensitivity of the WY procedure, and the instability of the random effects model.

The goal of the WY procedure is to estimate the adjusted p-value for outcome  $i$

$$\tilde{p}_i = Pr \left( \min_{1 \leq j \leq k} P_j \leq p_i \mid H_0^C \right).$$

where  $P_j$  is a random variable representing a  $p$ -value for outcome  $j$ , and lowercase  $p_i$  is the observed realization for outcome  $i$ . We assume a set of  $k$  tests with corresponding null hypotheses  $H_{0i}$  for  $i = 1, \dots, k$ . The complete null hypothesis is the setting where all the null hypotheses are true:  $H_0^C = \cap_{i=1}^k H_{0i} = \{\text{all } H_i \text{ are true}\}$ . In order to estimate this p-value, the p-value is calculated across  $B$  permutations

$$\tilde{p}_i = \frac{1}{B} \sum_{b=1}^B \mathbb{1}(\min_{1 \leq j \leq k} p_{b,j}^* \leq p_i).$$

In order for the procedure to be correct, we should be generating the p-values  $P_j$  from the the true distribution of p-values under the null hypothesis. We generally are testing outcomes that are truly significant, so the observed p-values are small. Thus, looking at these expressions, we are concerned about tail behavior; we are testing whether a small observed p-value is less than the minimum of of the generated p-values of our multiple outcomes. This combination of factors means that a small discrepancy in the tails between the generated null distribution and the true distribution could have a large impact on the adjusted p-value.

Given that we are relying on tail behavior, this explains why we need both a large number of permutations, and a large number of blocks/clusters. The large number of permutations is required because we are trying to estimate a rare event. With a significant outcome, it is rare that the observed p-value will be less than the minimum of p-values generated from the null distribution. The large number of blocks/clusters is required to accurately estimate the tail behavior for random effects models. With a small number of blocks/clusters, we may not properly estimate the spread of the distribution, resulting in a poor estimate of the tail behavior.