

# Heuristic Bias Correction - Conditional Poisson Distribution

April 6, 2019

## Contents

<b>1</b>	<b>Description</b>	<b>1</b>
<b>2</b>	<b>Bias Function</b>	<b>5</b>
<b>3</b>	<b>Alternative bias-corrected estimators</b>	<b>8</b>
3.1	Bias in % . . . . .	8
3.1.1	$\hat{\psi}$ . . . . .	8
3.2	CV in percent . . . . .	14
<b>4</b>	<b>Ratio of bias to sd</b>	<b>27</b>
4.1	$\hat{\lambda}$ . . . . .	27
<b>5</b>	<b>Histograms</b>	<b>33</b>
<b>6</b>	<b>Boxplots</b>	<b>58</b>

## 1 Description

The MLE is only meaningful for non-degenerate data. However, the bias correction does not take this fact properly into account.

The probability of obtaining data with  $\sum_{k=1}^n N_k = N$  is given by

$$q_1 := \left( \sum_{j=1}^n Q_{\mathbf{e}_j} \right)^N, \quad (1)$$

where

$$Q_{\mathbf{e}_j} = \frac{e^{\lambda p_j} - 1}{e^{\lambda} - 1}. \quad (2)$$

The probability that only one lineage is observed in the data is given by

$$\sum_{j=1}^n Q_{\mathbf{e}_j}^N.$$

In this case there are no signs of super-infections. Hence, the MLE is not uniquely defined, rather all admissible values of  $\lambda$  are equally likely given that  $p_j = 1$  if only lineage  $A_j$  is observed. The probability of obtaining data with no signs of super-infection but at least two different lineages present,  $\sum_{k=1}^n N_k = N$  and  $N_k \neq N$  for all  $k = 1, 2, \dots, n$ , is given by

$$q_2 := \left( \sum_{j=1}^n Q_{\mathbf{e}_j} \right)^N - \sum_{j=1}^n Q_{\mathbf{e}_j}^N. \quad (3)$$

In this case the MLE would formally be  $\hat{\lambda} = 0$ . The probability of obtaining data with  $N_k = N$  for at least one  $k$  is given by

$$q_3 := \frac{1}{(1 - e^{-\lambda})^N} \left( 1 - \prod_{j=1}^n (1 - (1 - e^{-\lambda p_j})^N) \right) \quad (4)$$

in which case no MLE exists. In summary the probability of having degenerate data becomes

$$q := q_3 + q_2 \quad (5)$$

These facts are not properly addressed by the bias correction. A regular dataset  $\mathbf{X}$ , where  $\mathbf{X}$  is a 0-1 matrix of size  $N \times n$ , is a dataset where  $N_k = \sum_{j=1}^n \mathbf{X}_{kj} \neq N$  for  $k = 1, 2, \dots, n$  and  $\sum_{k=1}^n N_k = \sum_{k=1}^n \sum_{j=1}^n \mathbf{X}_{kj} > N$ . Let  $\mathcal{X}$  be the set of all regular datasets. By iterated expectation we have

$$\mathbb{E}\{\hat{\lambda}^{(\text{bc})}\} = \mathbb{E}\{\mathbb{E}[\hat{\lambda}^{(\text{bc})} | \mathbb{1}_{\mathcal{X}}(\mathbf{X})]\} \quad (6a)$$

$$= \Pr\{\mathbb{1}_{\mathcal{X}}(\mathbf{X}) = 1\} \mathbb{E}\{\hat{\lambda}^{(\text{bc})} | \mathbb{1}_{\mathcal{X}}(\mathbf{X}) = 1\} + \Pr\{\mathbb{1}_{\mathcal{X}}(\mathbf{X}) = 0\} \mathbb{E}\{\hat{\lambda}^{(\text{bc})} | \mathbb{1}_{\mathcal{X}}(\mathbf{X}) = 0\} \quad (6b)$$

$$\approx \Pr\{\mathbb{1}_{\mathcal{X}}(\mathbf{X}) = 1\} \mathbb{E}\{\hat{\lambda}^{(\text{bc})} | \mathbb{1}_{\mathcal{X}}(\mathbf{X}) = 1\} \quad (6c)$$

where

$$\Pr\{\mathbb{1}_{\mathcal{X}}(\mathbf{X}) = 1\} = 1 - q \quad (7)$$

is the probability of regular data. Asymptotically equality holds in 6c. However, if  $N$  is small and  $\lambda$  moderate and the lineage frequencies are highly skewed, second term in 6c contributes substantially to the sum. In moderate settings, a better bias-corrected estimate would be

$$\tilde{\lambda}^{(\text{bc})} = (1 - q)\hat{\lambda}^{(\text{bc})} \quad (8)$$

where  $q_3$  and  $q_2$  are evaluated at MLE. If the bias function is non-constant, preferably the analytic bias expression should be evaluated at BCMLE (MacKinnon and Smith-1998). Consequently a better estimate for  $q$  is derived at the bias-corrected estimate  $\hat{\lambda}^{(\text{bc})}$ .

**Remark 1.** *A better bias-corrected estimate of MOI parameter is*

$$\tilde{\lambda}^{(bc)} = (1 - q)\hat{\lambda}^{(bc)}, \quad (9)$$

*where  $q$  is the probability of having pathological data. Similarly, we can obtain better estimations for lineage frequencies*

$$\hat{p}_k^{bc} = (1 - q)\hat{p}_k^{bc}. \quad (10)$$

Let  $r = \frac{1-q_3}{1-q_1}$  an alternative estimate is introduced in the next remark.

**Remark 2.** *An alternative bias-corrected estimate is*

$$\tilde{\lambda}^{(bc)} = \hat{\lambda} - \text{bias}(\hat{\lambda})r, \quad (11)$$

*where  $q_1$ ,  $q_3$  and  $\text{bias } \hat{\lambda}$  are evaluated at  $\hat{\boldsymbol{\theta}}^{(bc)}$ .*



## 2 Bias Function

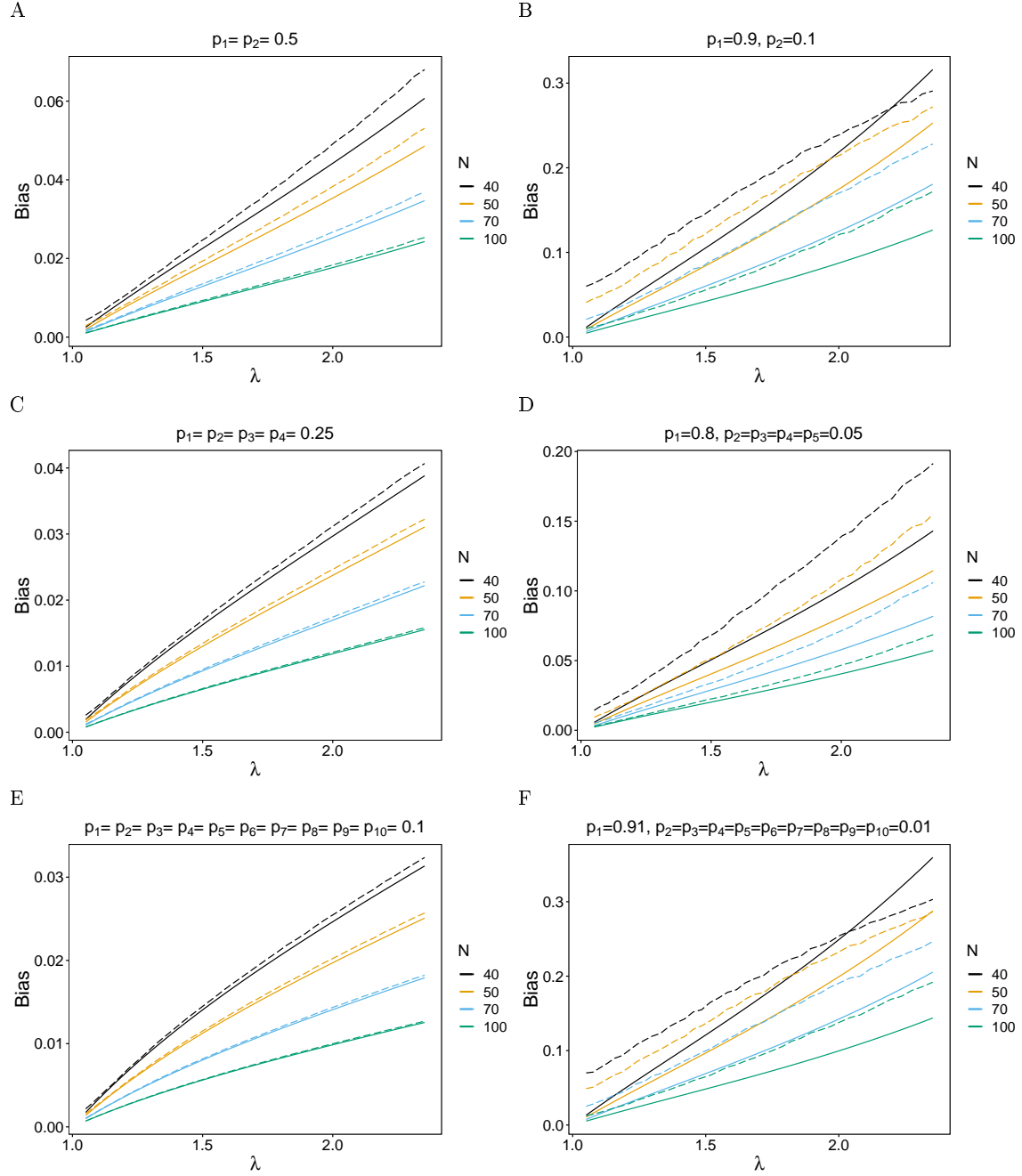


Figure 1: **Bias function.** The figure shows the behavior of average approximation of theoretical bias from  $S = 10,000$  simulated datasets (solid lines) versus the theoretical bias at the true parameter (long-dashed lines). Different sample sizes are specified with different colors.

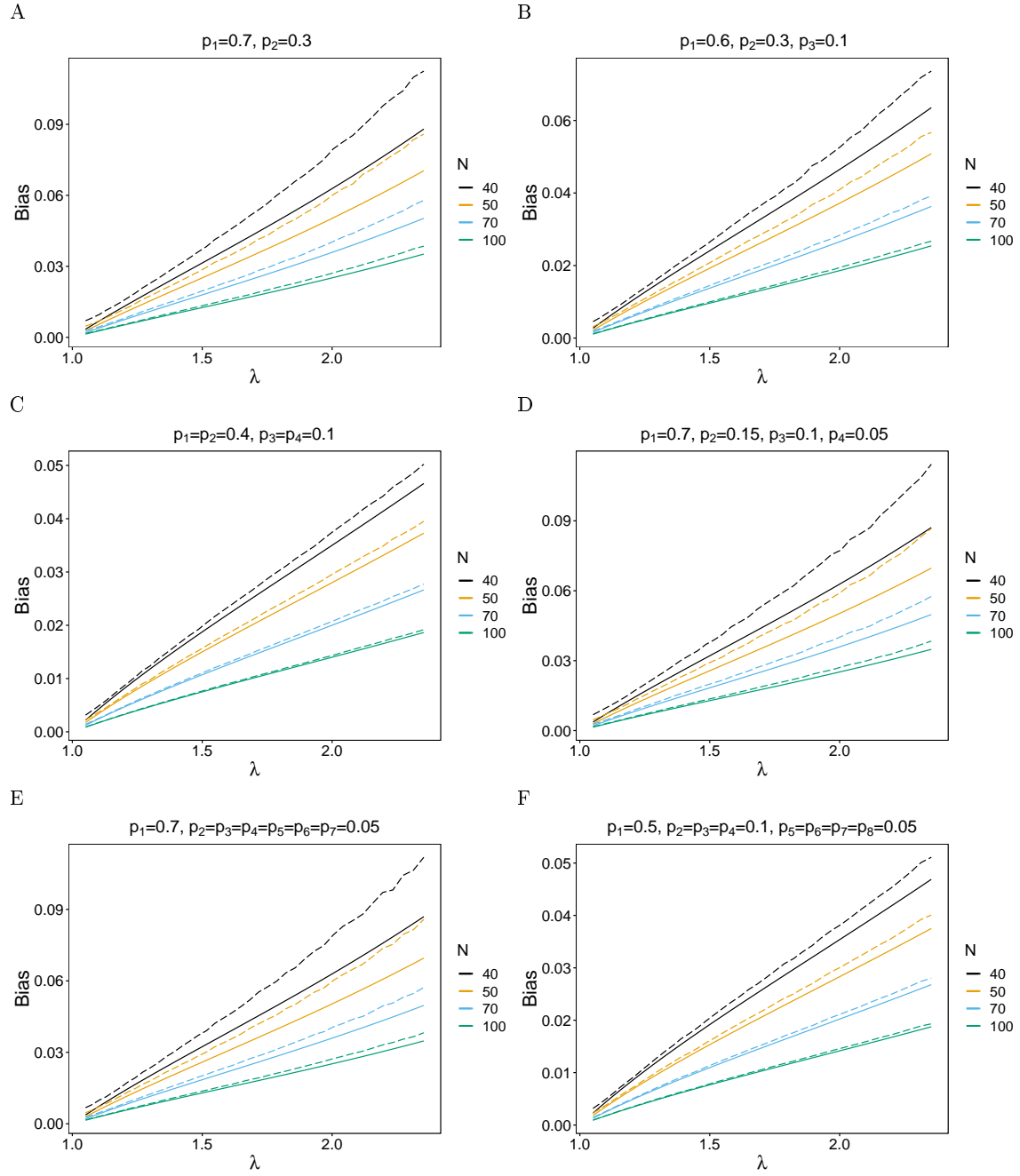


Figure 2: **Bias function.** Same as Figure 1



### 3 Alternative bias-corrected estimators

#### 3.1 Bias in %

##### 3.1.1 $\hat{\psi}$

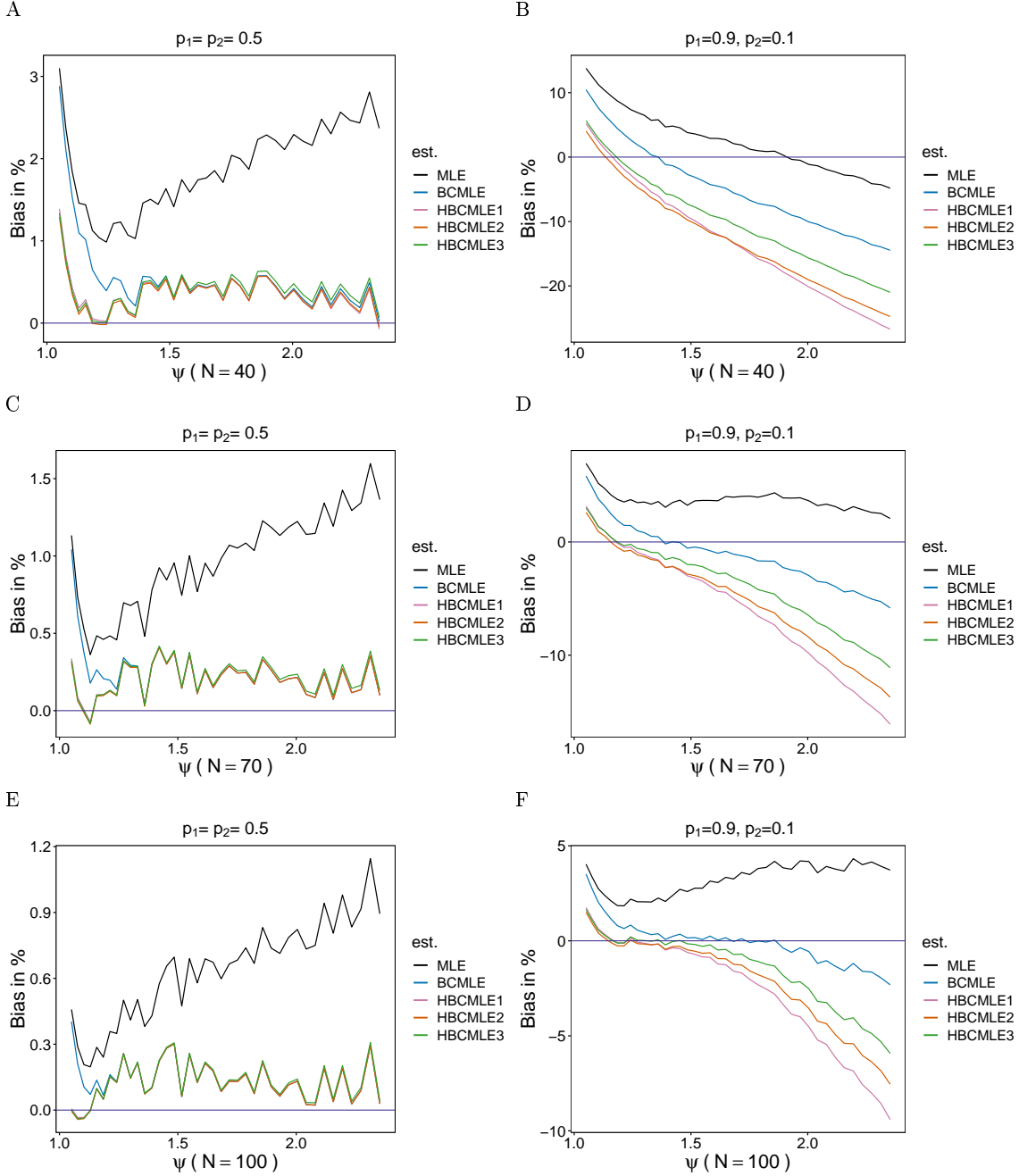


Figure 3: **Bias of alternative estimators in %**. The figure shows the bias in % of different estimators of  $\psi = \frac{\lambda}{1-e^{-\lambda}}$  (average MOI). We have  $\text{HBCMLE1} = (\hat{\lambda} - \text{bias}(\hat{\lambda})|_{\lambda=\hat{\lambda}})(1 - q_1 - q_3 + q_2)$  where  $q_1, q_2$  and  $q_3$  are evaluated at MLE. Further,  $\text{HBCMLE2} = \hat{\lambda} - \text{bias}(\hat{\lambda})|_{\lambda=\hat{\lambda} \frac{1-q_3}{1-q_1}}$  and  $\text{HBCMLE3} = (\hat{\lambda} - \text{bias}(\hat{\lambda})|_{\lambda=\hat{\lambda}})(1 - q_1 - q_3 + q_2)$  and  $\text{HBCMLE4} = \hat{\lambda} * (1 - q_3 - q_1 + q_2) - \text{bias}(\hat{\lambda})|_{\lambda=\hat{\lambda}}$  and  $\text{HBCMLE5} = (\hat{\lambda} - \text{bias}(\hat{\lambda})|_{\lambda=\hat{\lambda}^{(b_c)}}) * (1 - q_3 - q_1 + q_2)$  where  $q_1, q_2$  and  $q_3$  are evaluated at BCMLE. The transformation  $\frac{\lambda}{1-e^{-\lambda}}$  is applied afterwards.



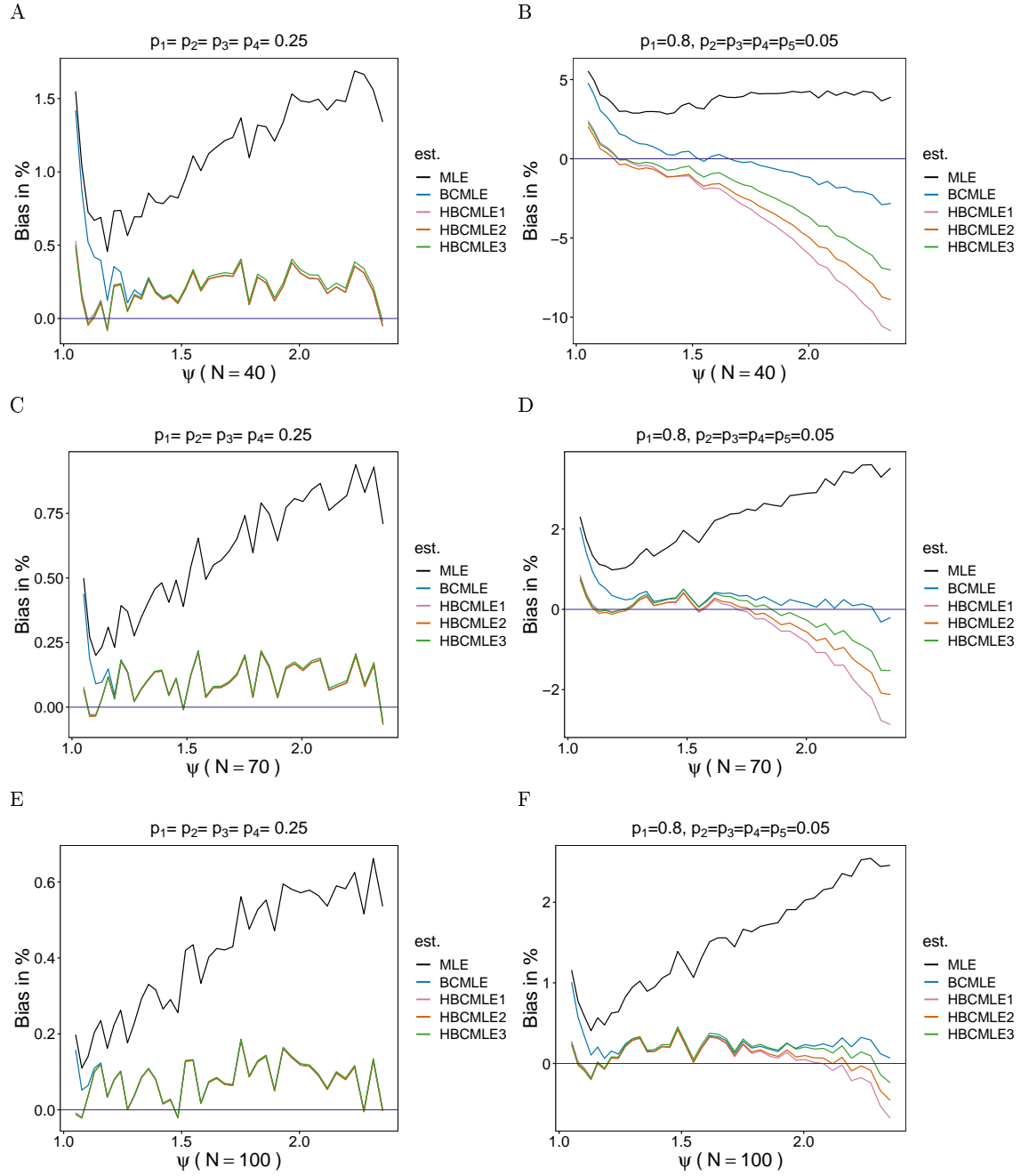


Figure 4: Same as Figure 3 but for different lineage-frequency distributions.

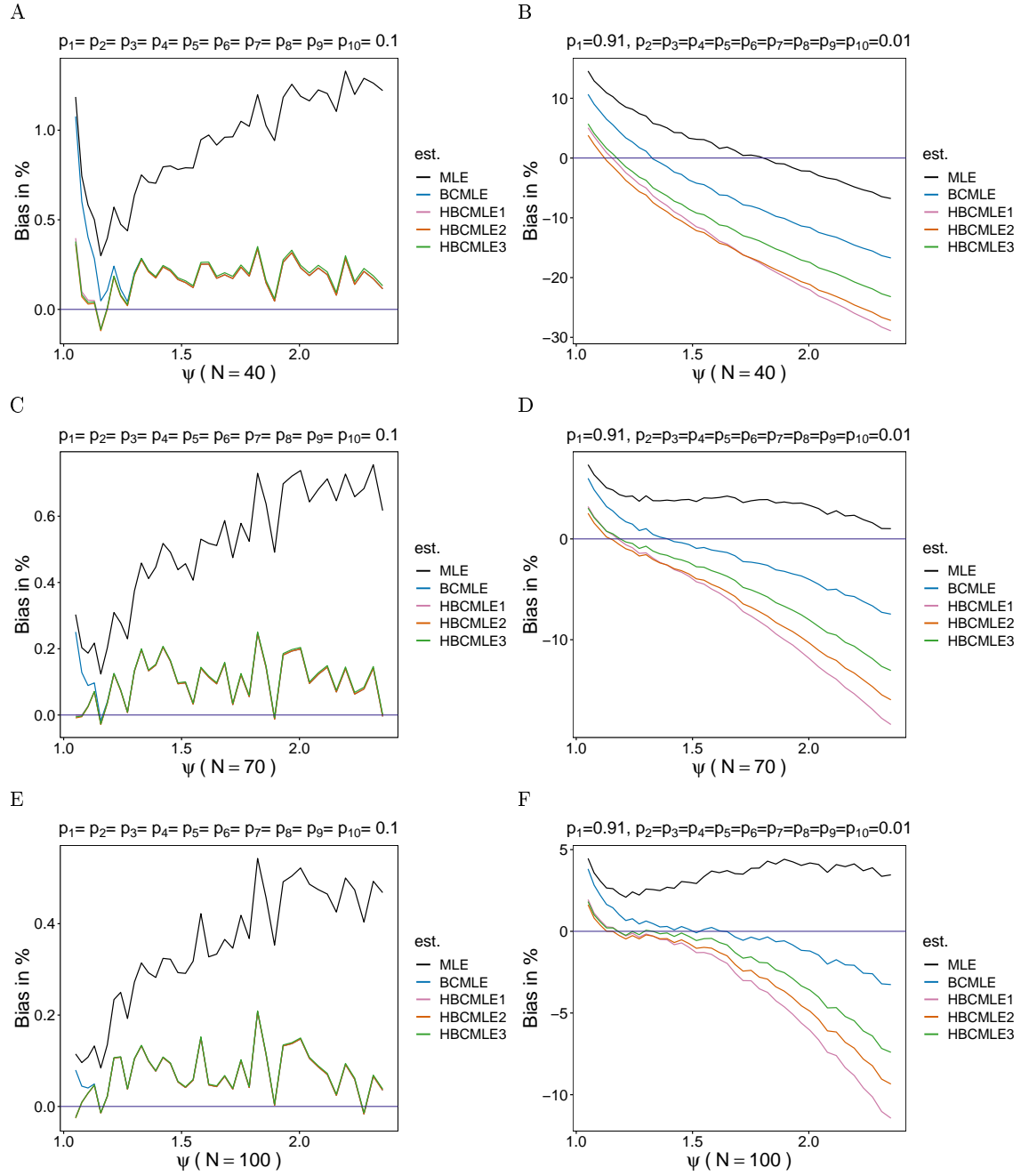


Figure 5: Same as Figure 3 but for different lineage-frequency distributions.

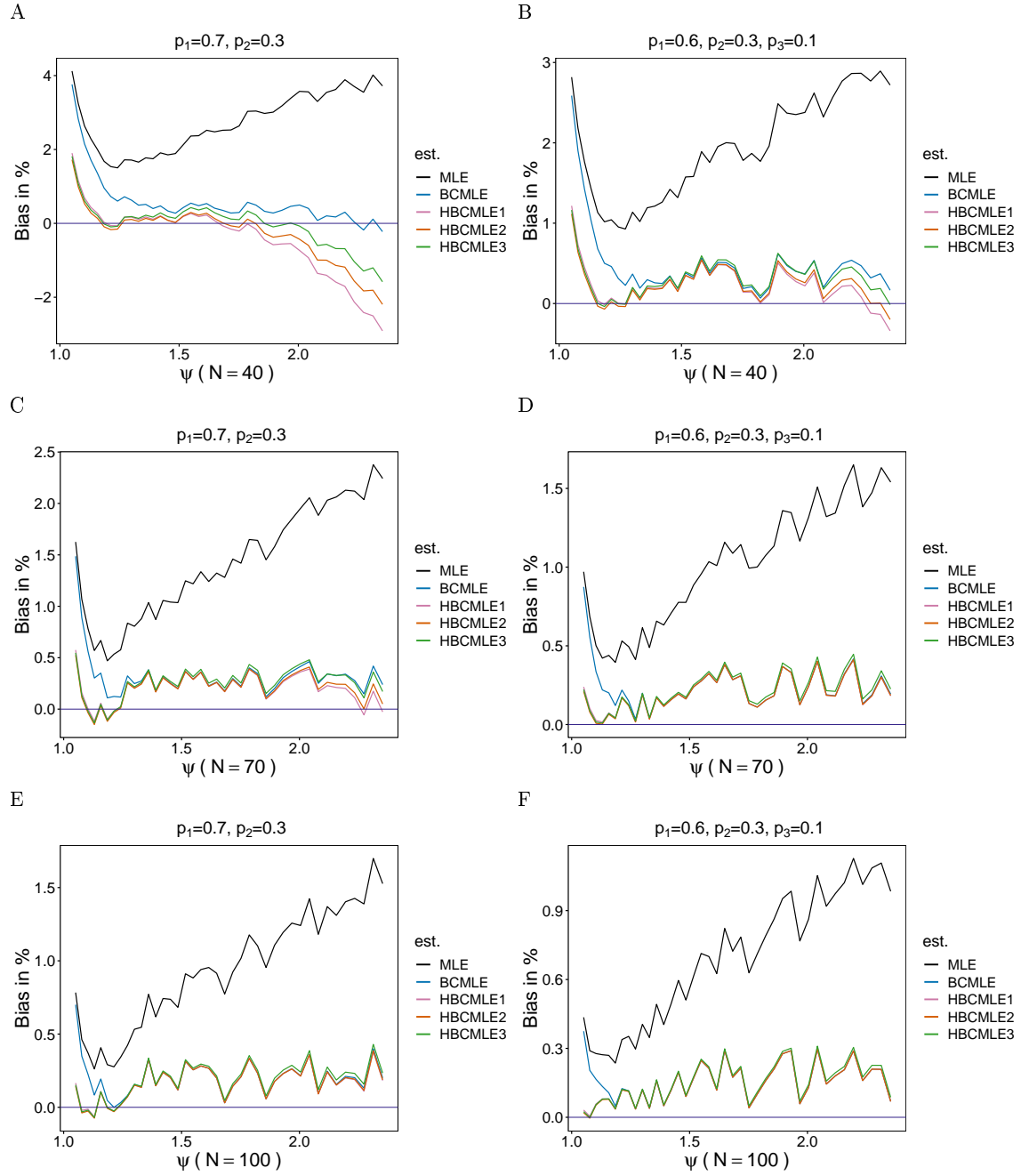


Figure 6: **Bias of alternative estimators in %.** Similar to Figure 3.

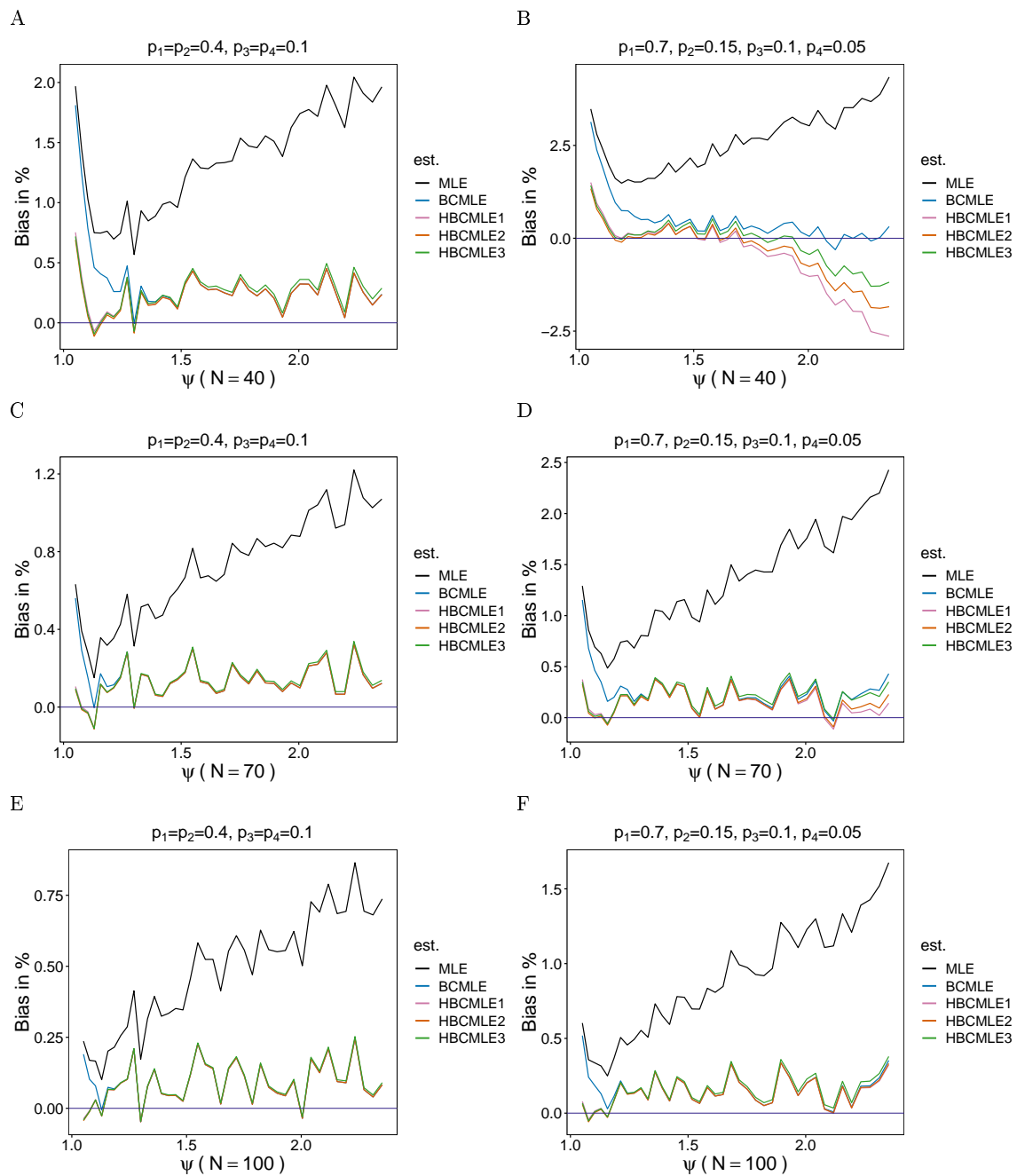


Figure 7: Same as Figure 3 but for different lineage-frequency distributions.

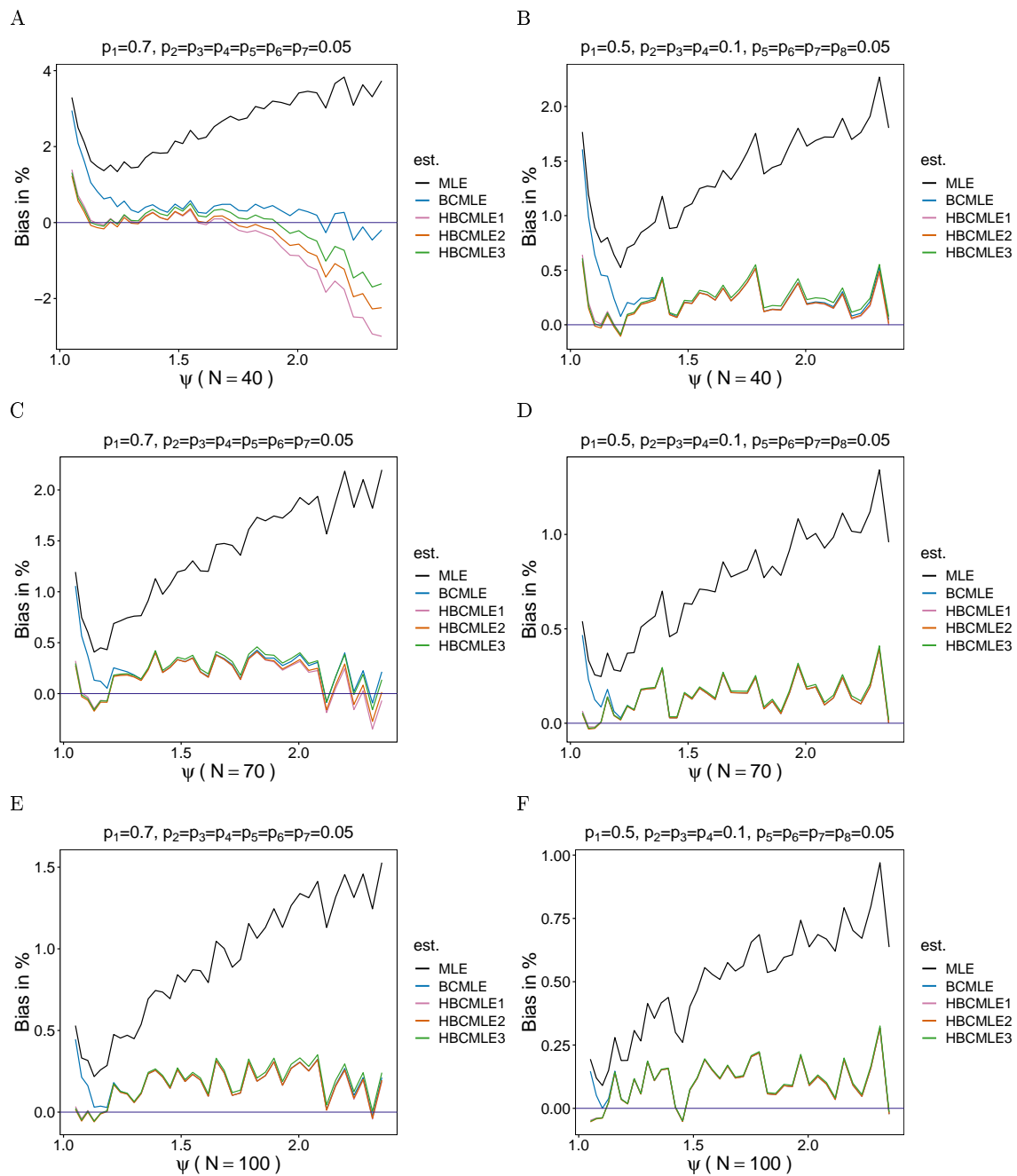


Figure 8: Same as Figure 3 but for different lineage-frequency distributions.

### 3.2 CV in percent

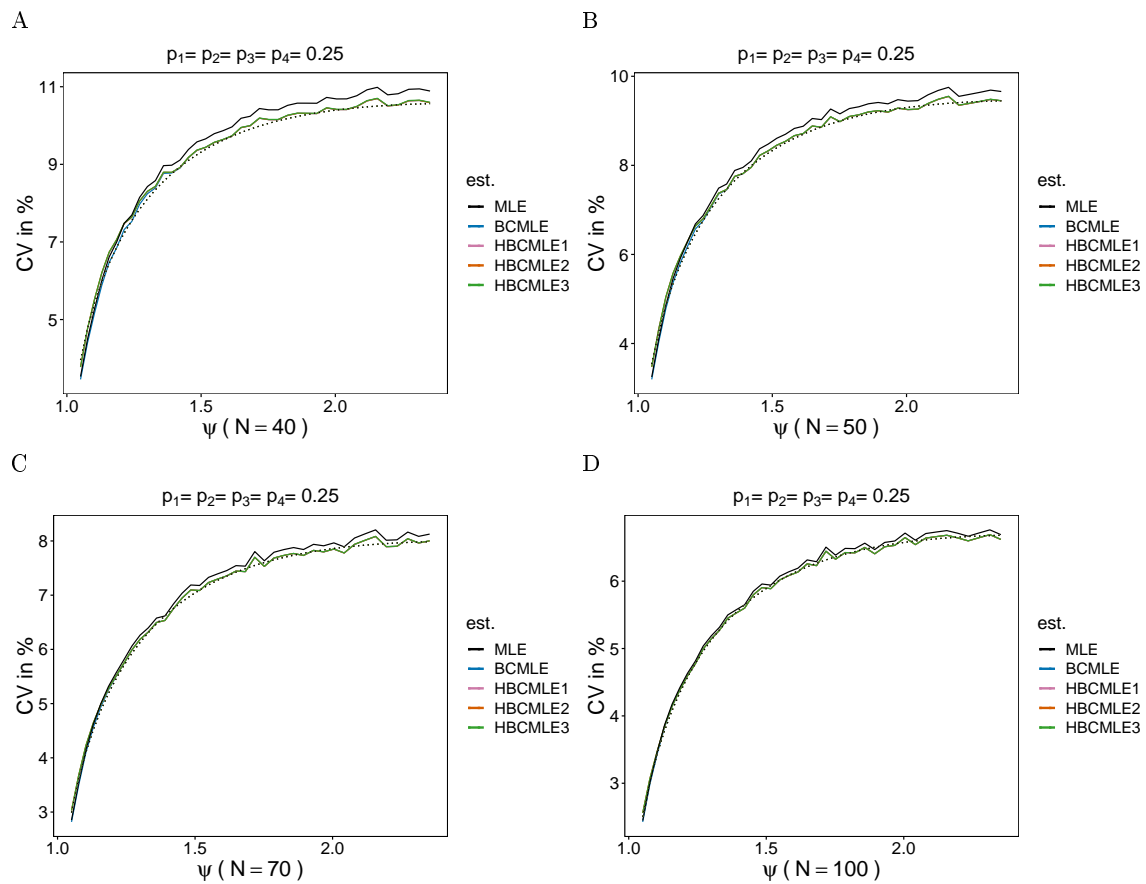
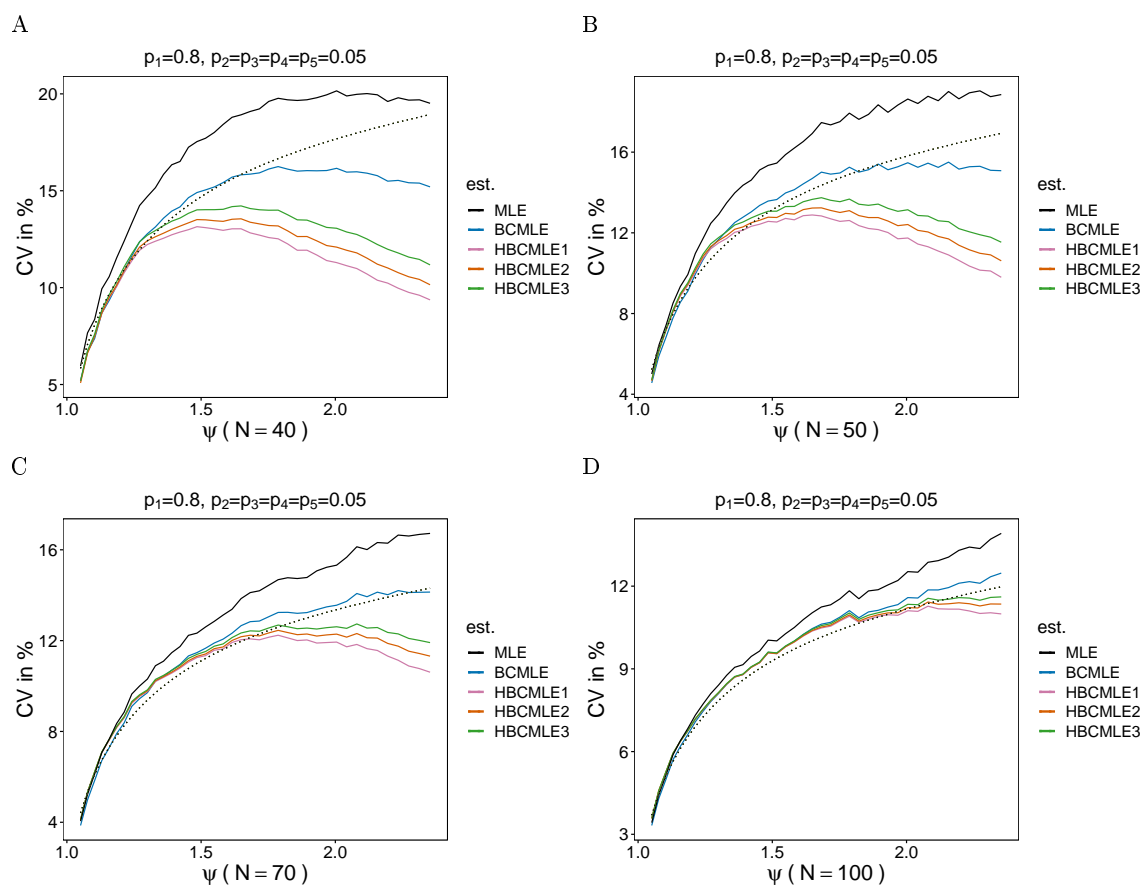
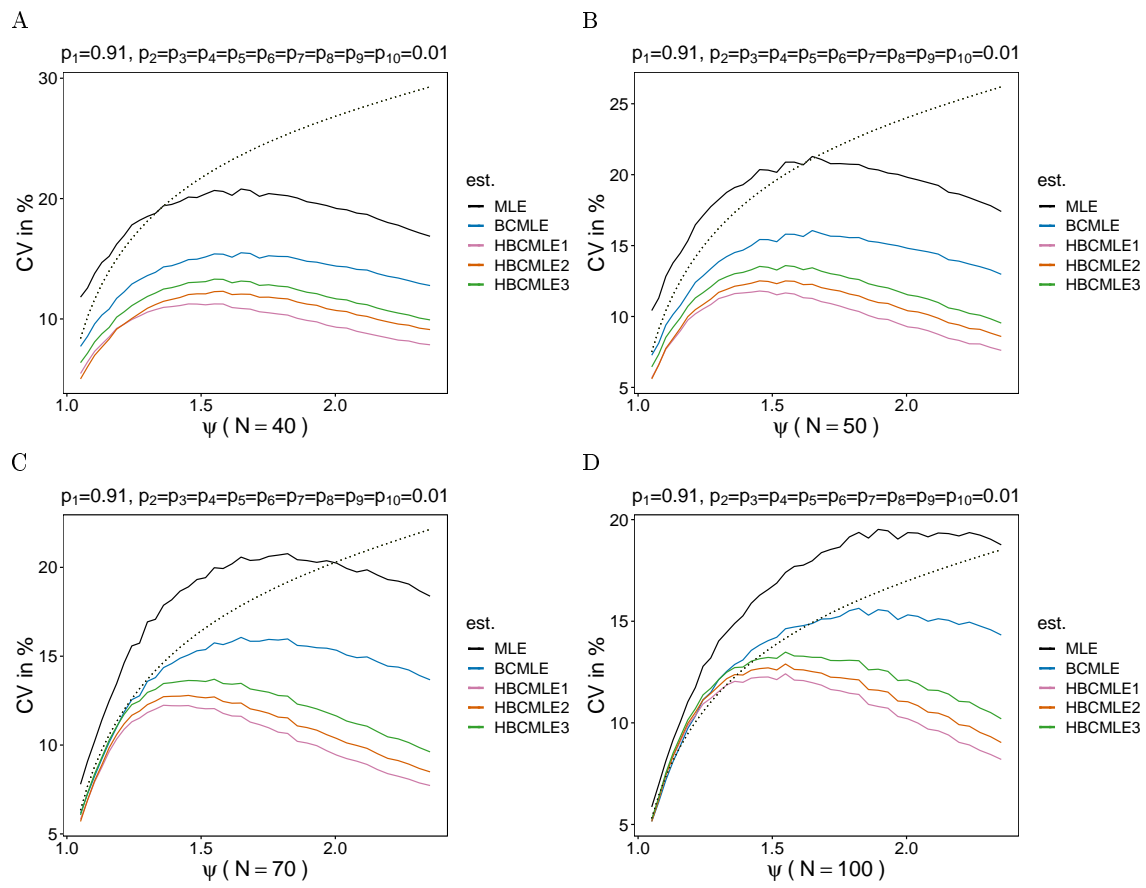


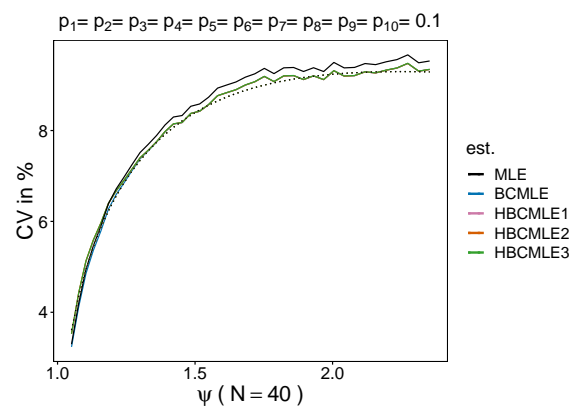
Figure 9: CV of alternative estimators in %. Same as Figure ?? but for CV



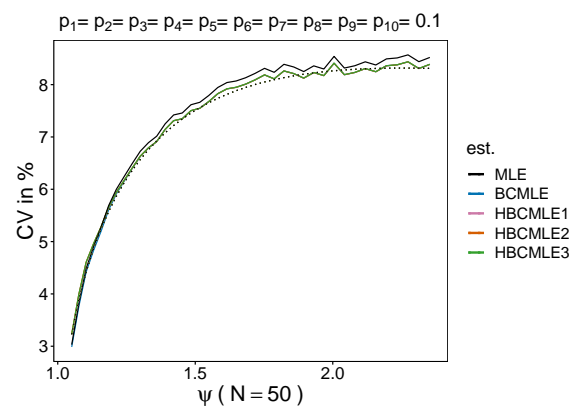




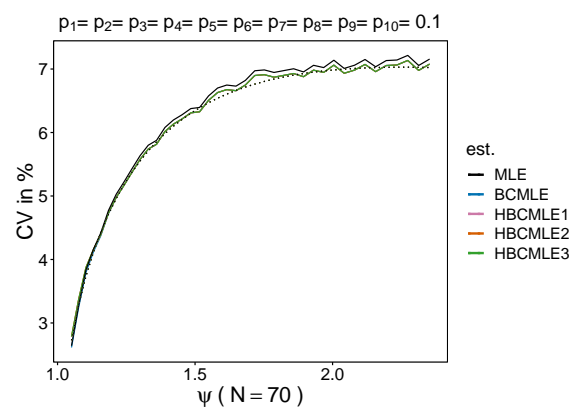
A



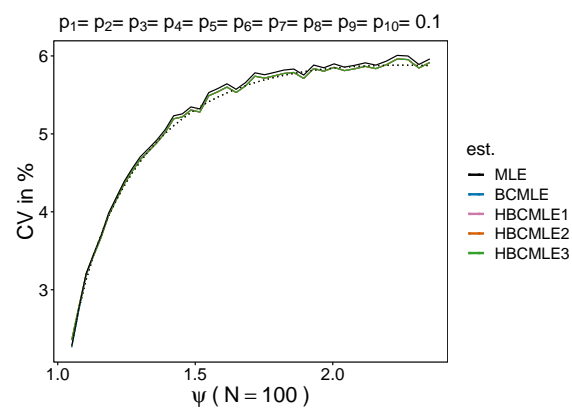
B



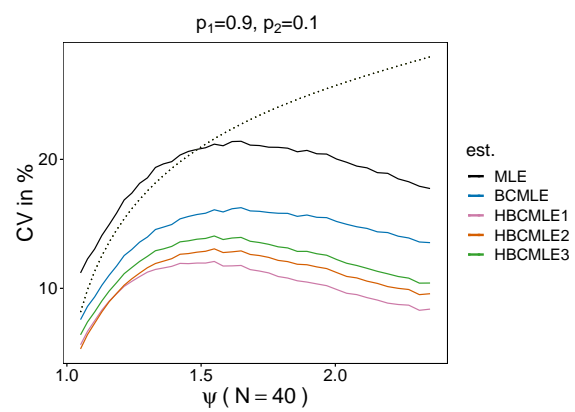
C



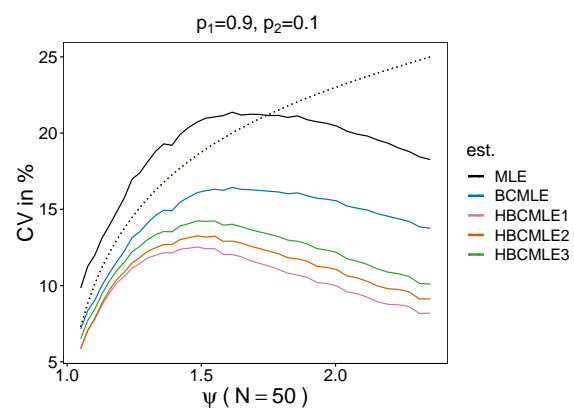
D



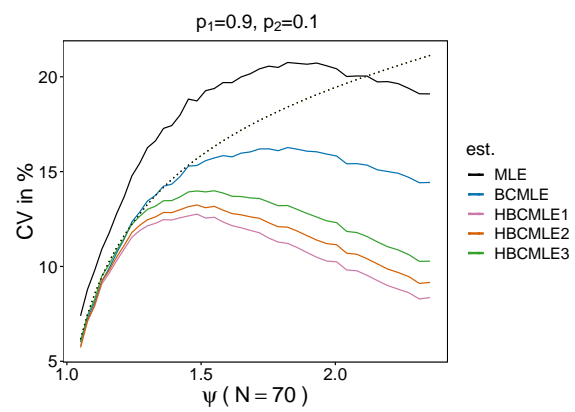
A



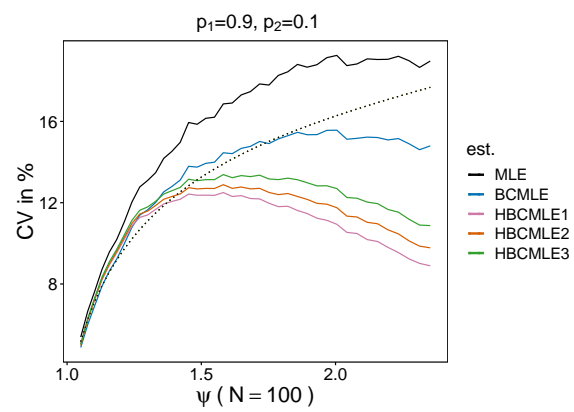
B

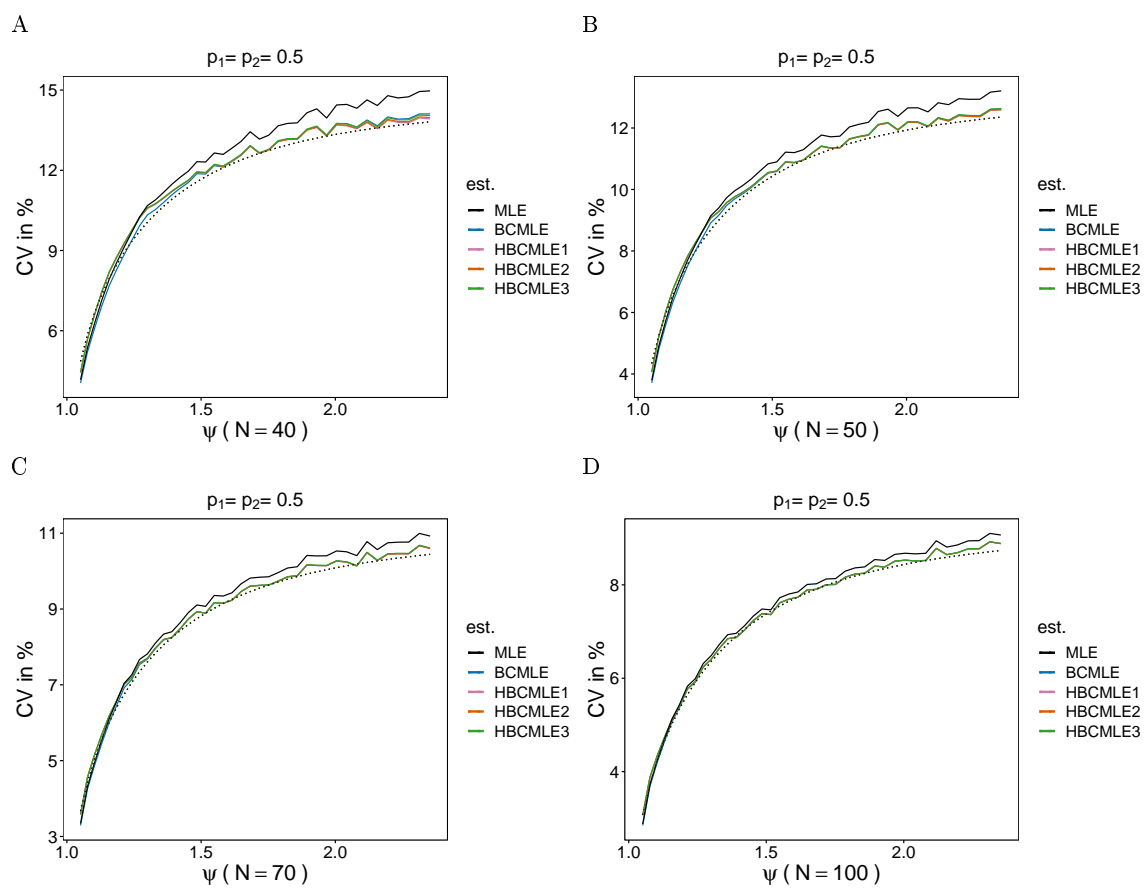


C

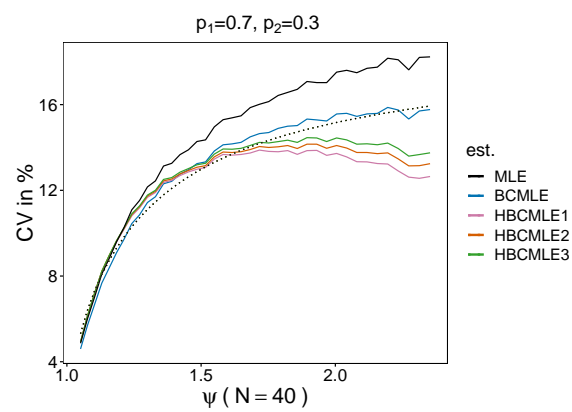


D

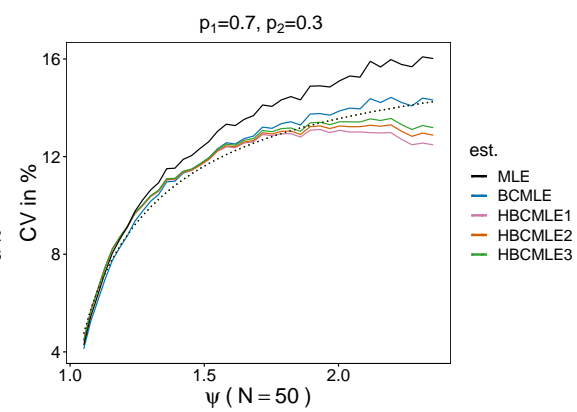




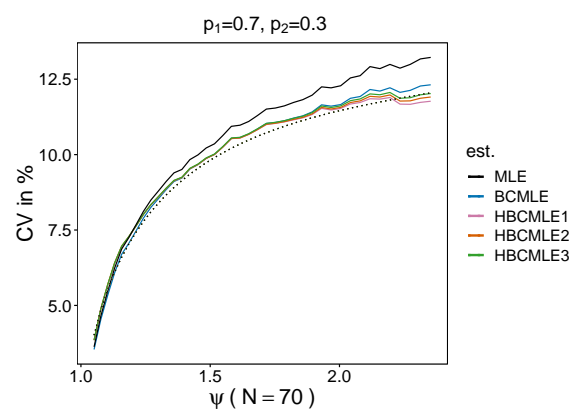
A



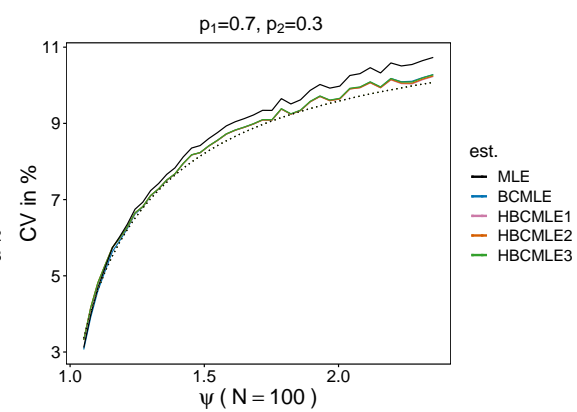
B

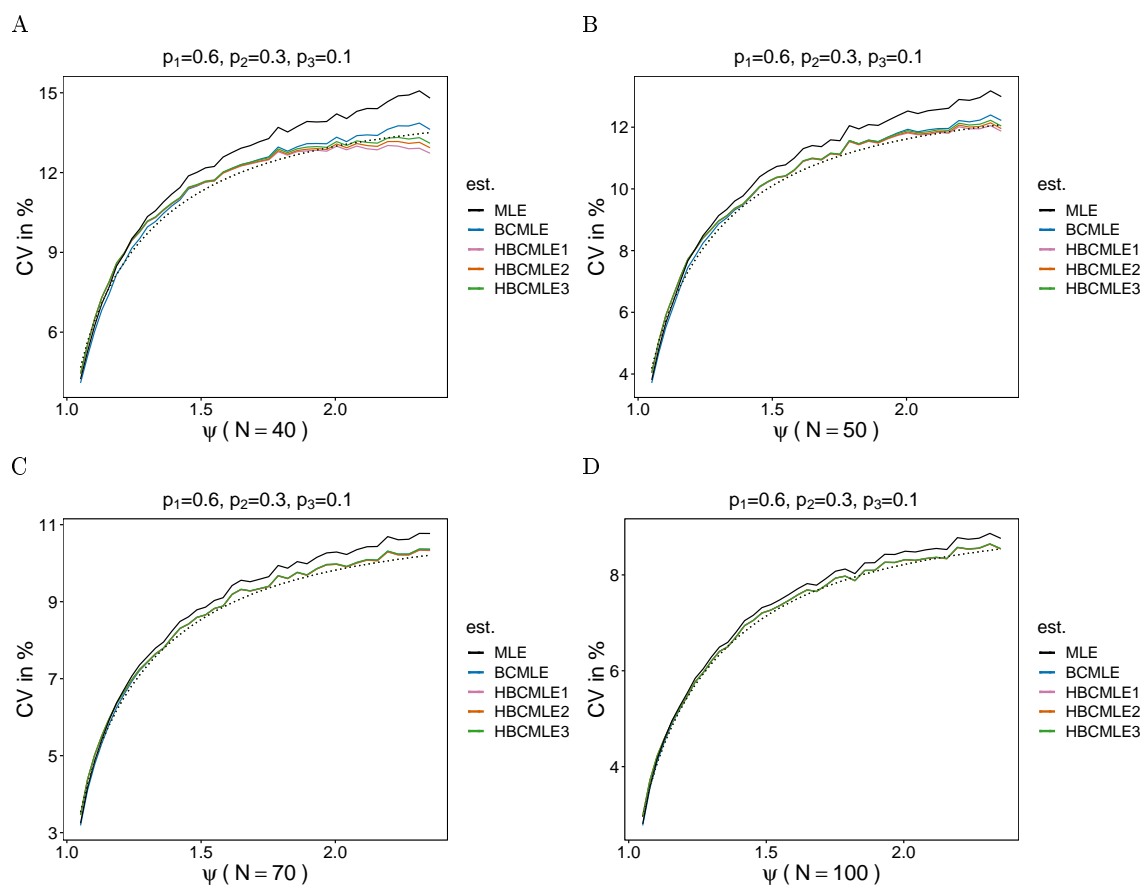


C

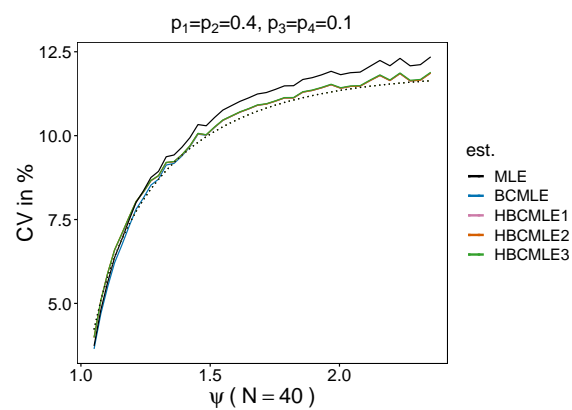


D

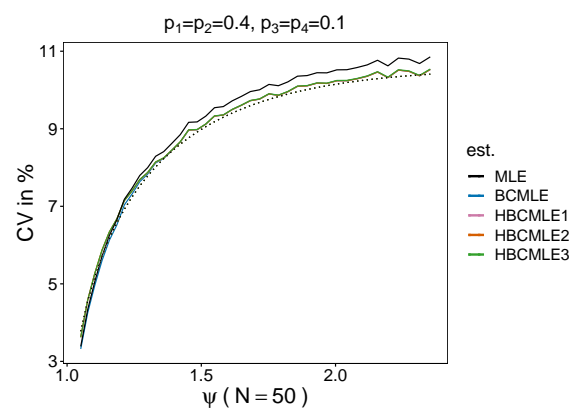




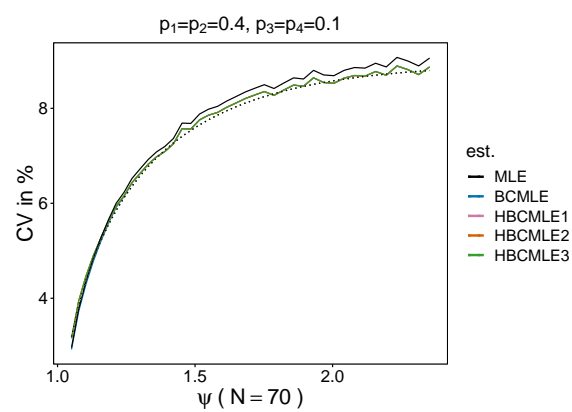
A



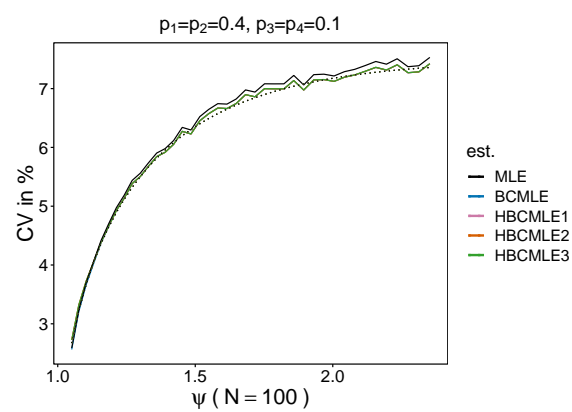
B

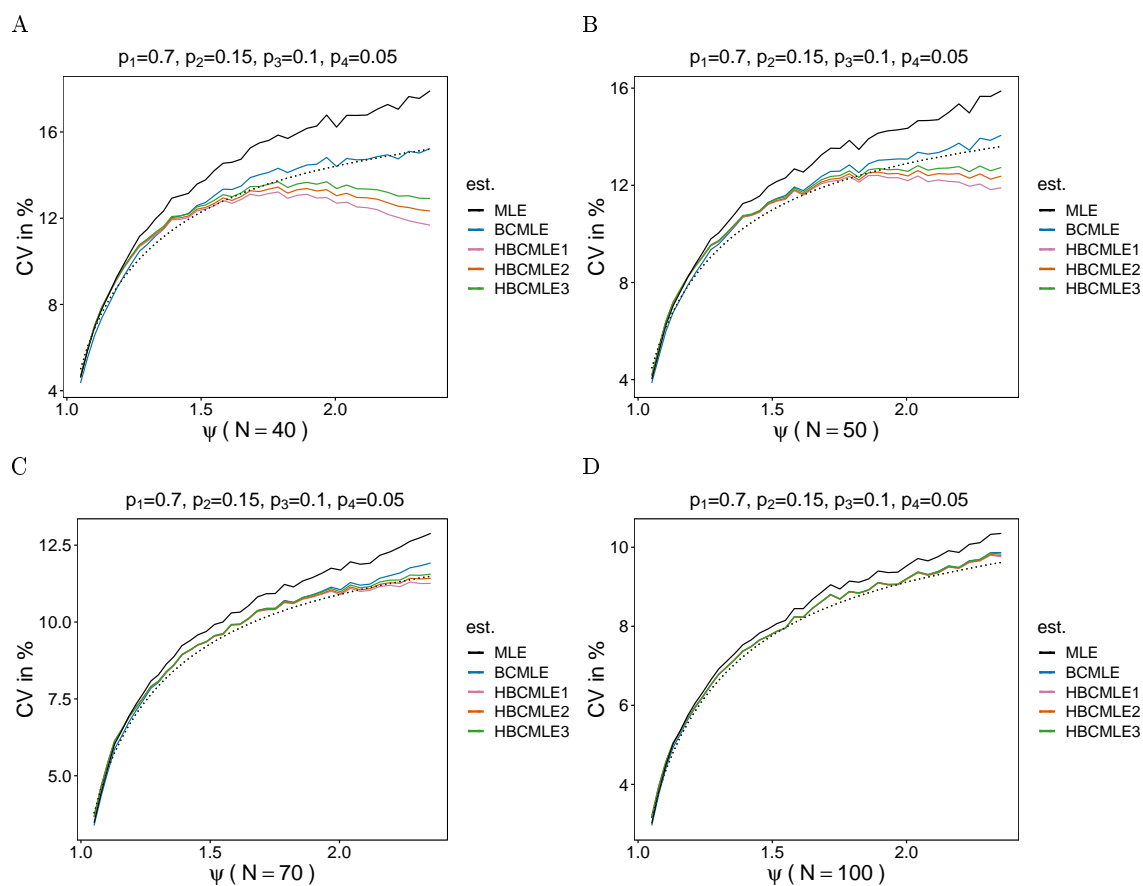


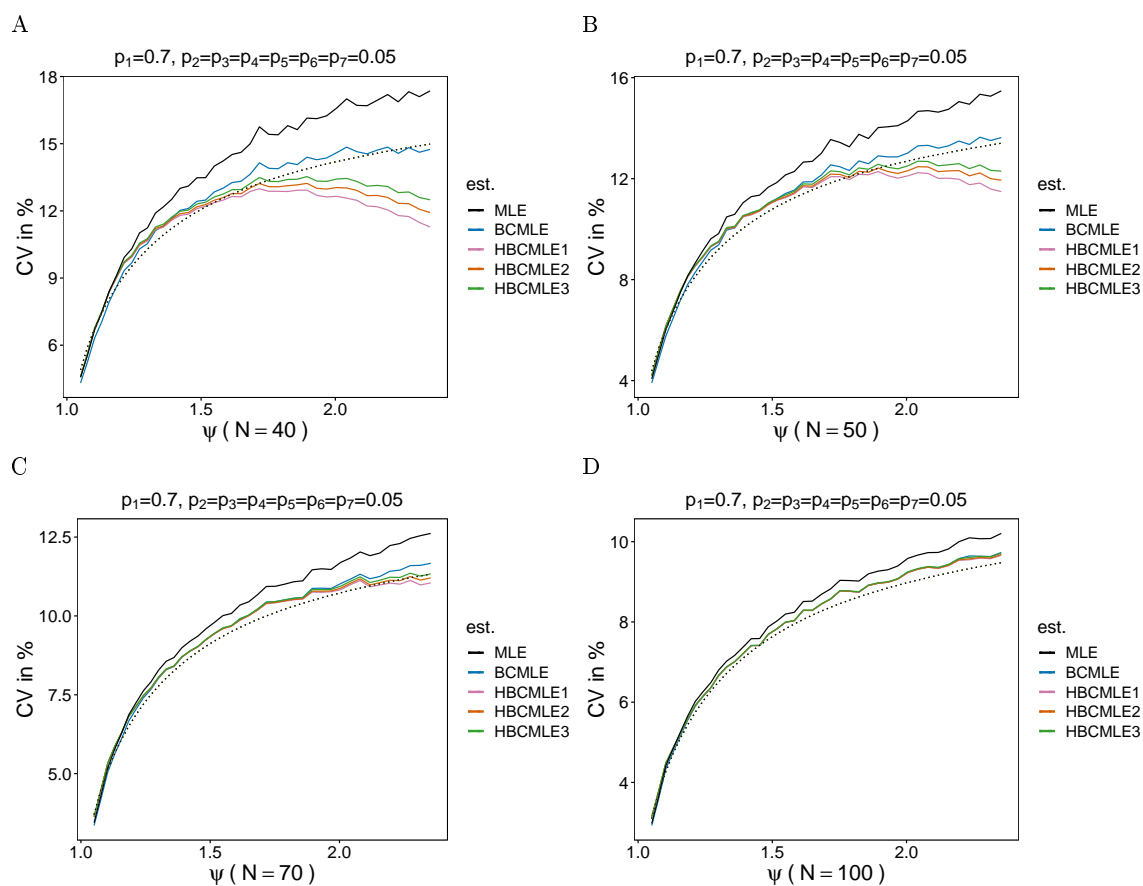
C



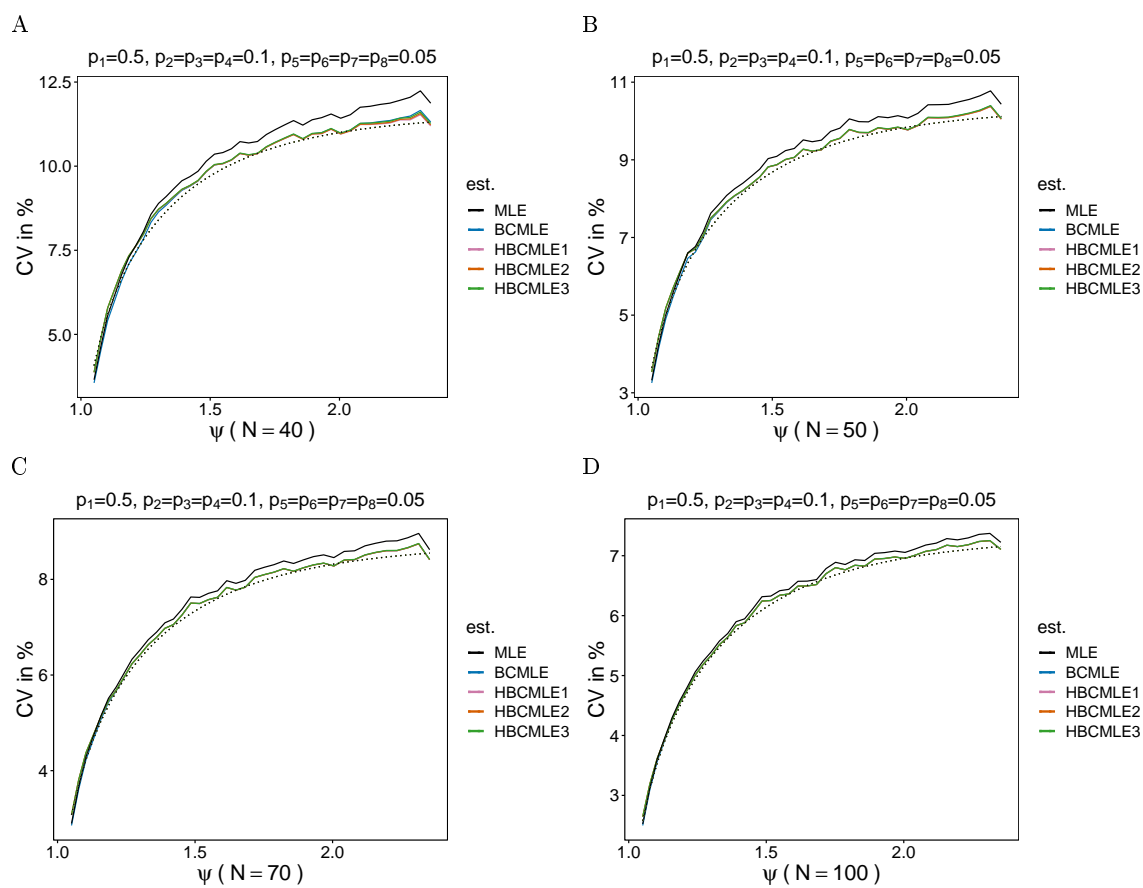
D













## 4 Ratio of bias to sd

### 4.1 $\hat{\lambda}$

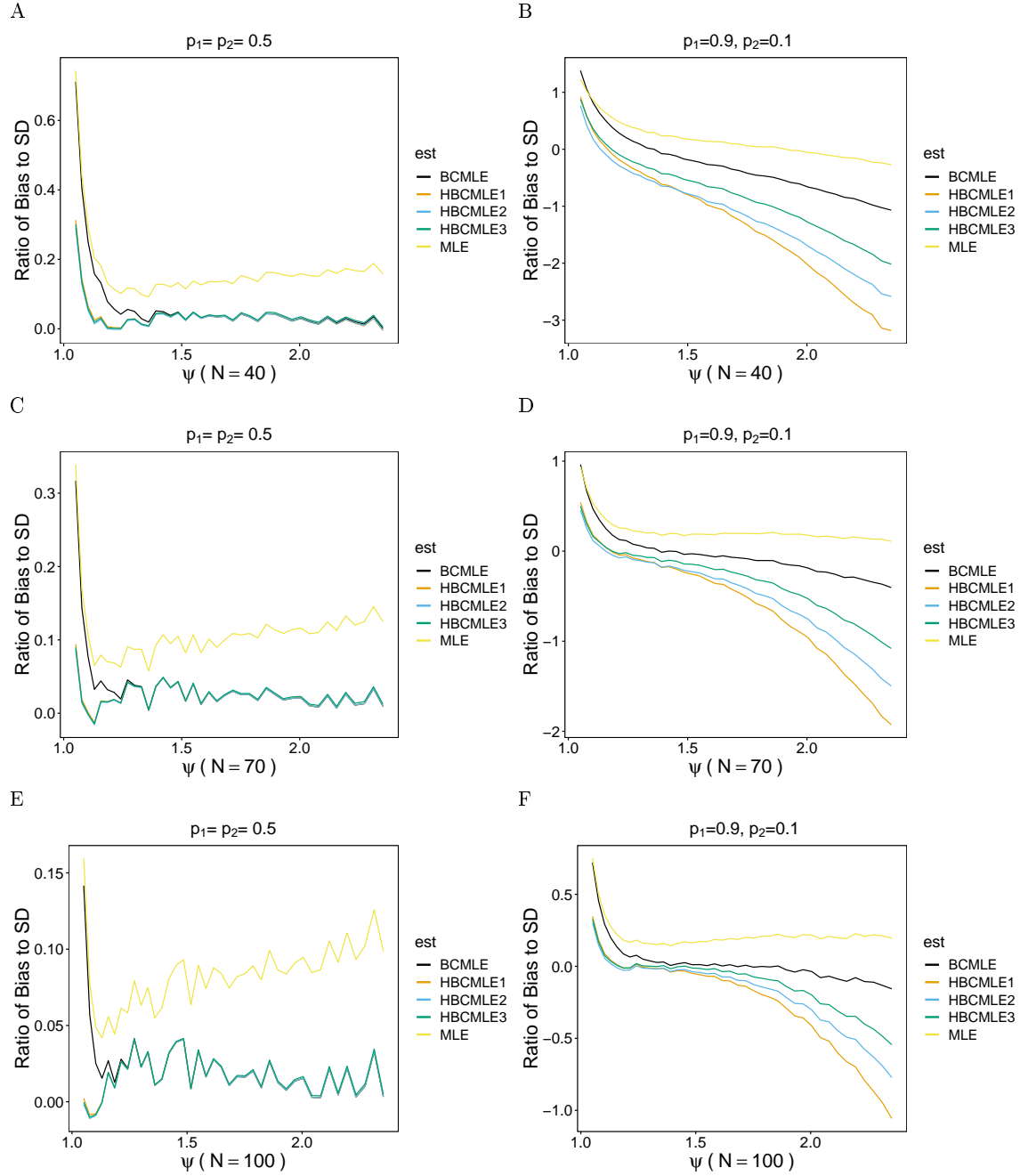


Figure 21: **Ratio of MLE to BCMLE.** The figure shows ratio of MLE to BCMLE. Different colors correspond to different ratios. The dotted and dashed lines correspond to true  $q$  5 and  $r$  2, respectively.

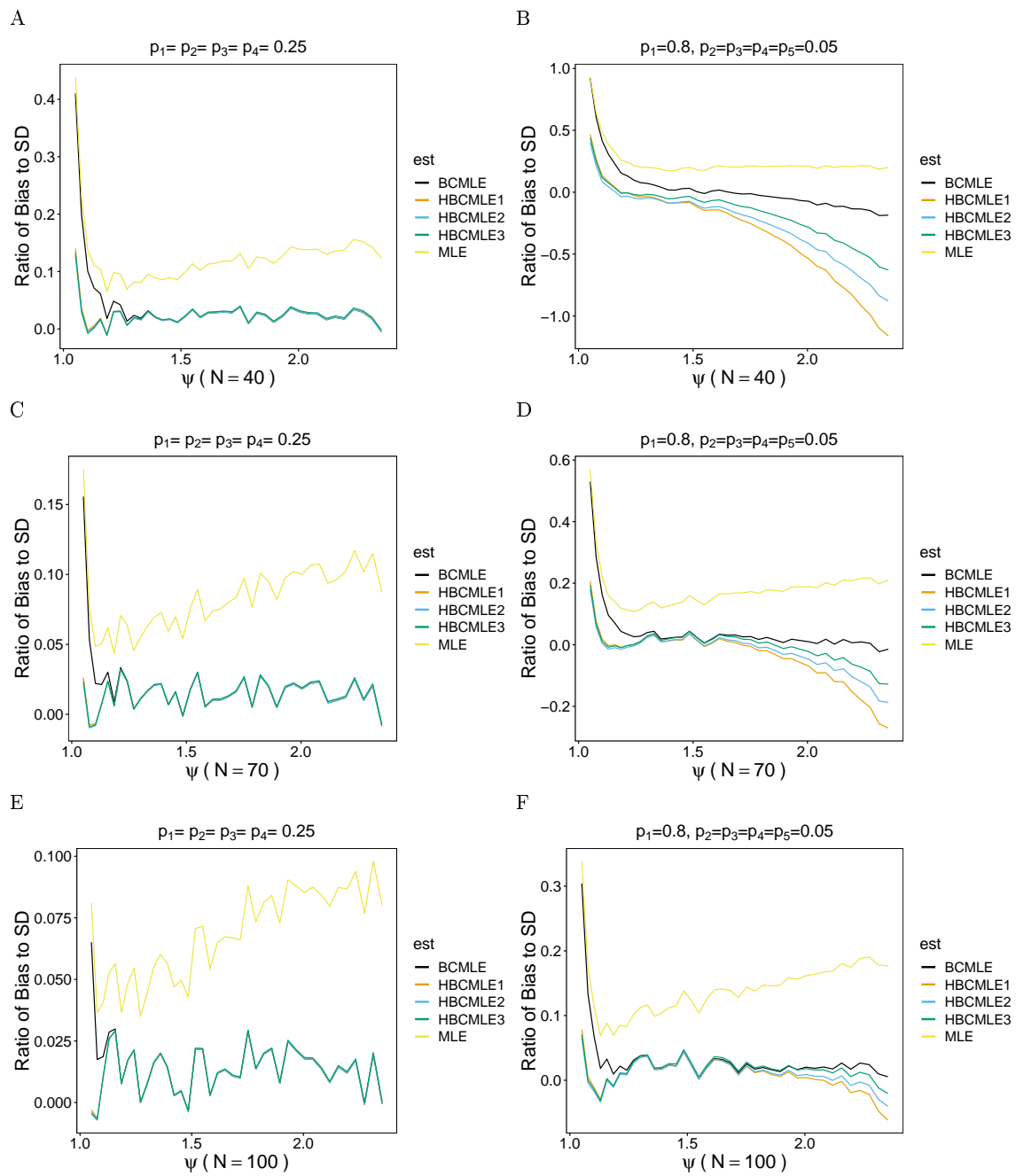


Figure 22: Same as Figure 21 but for different lineage-frequency distributions.

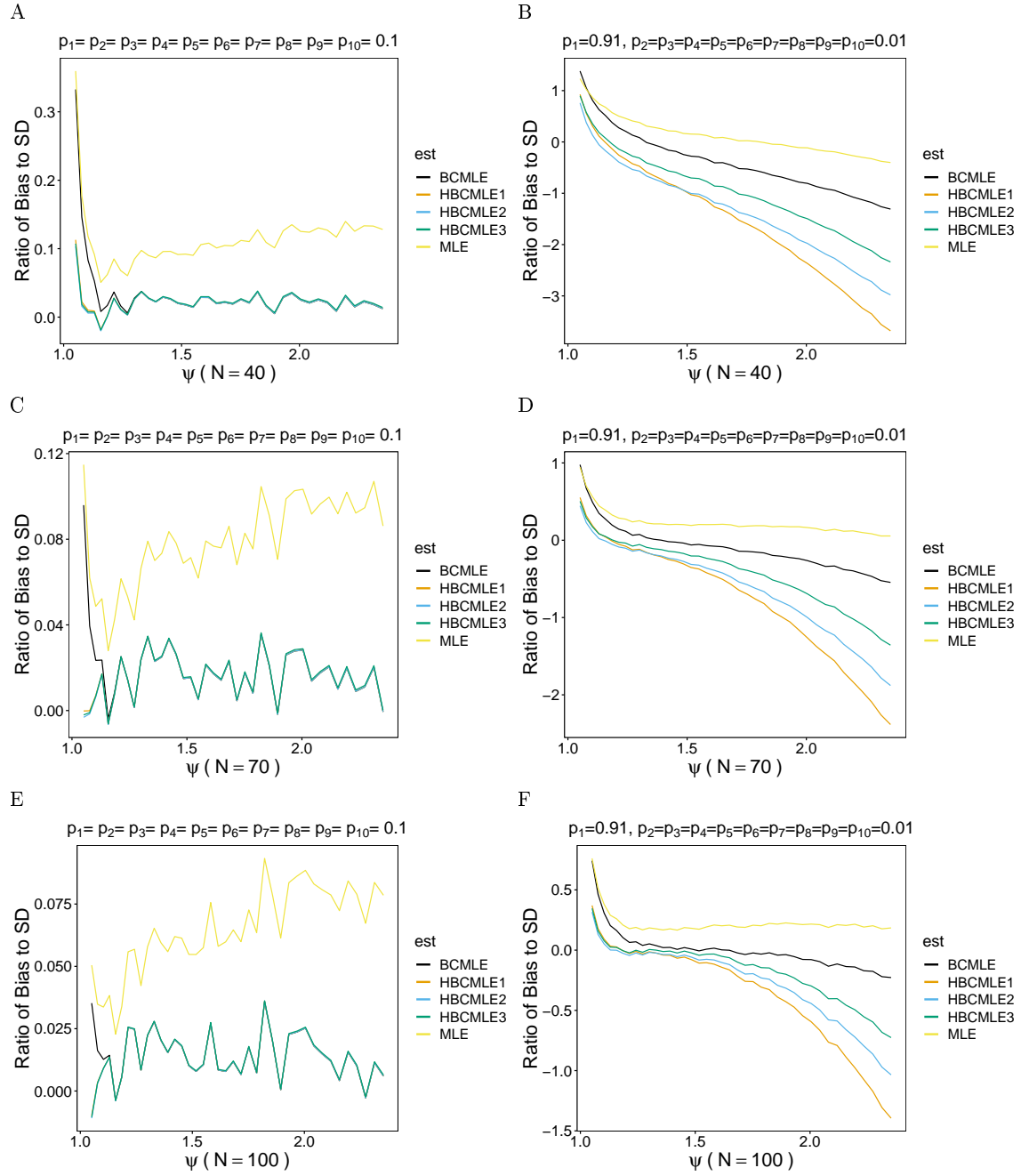


Figure 23: Same as Figure 21 but for different lineage-frequency distributions.

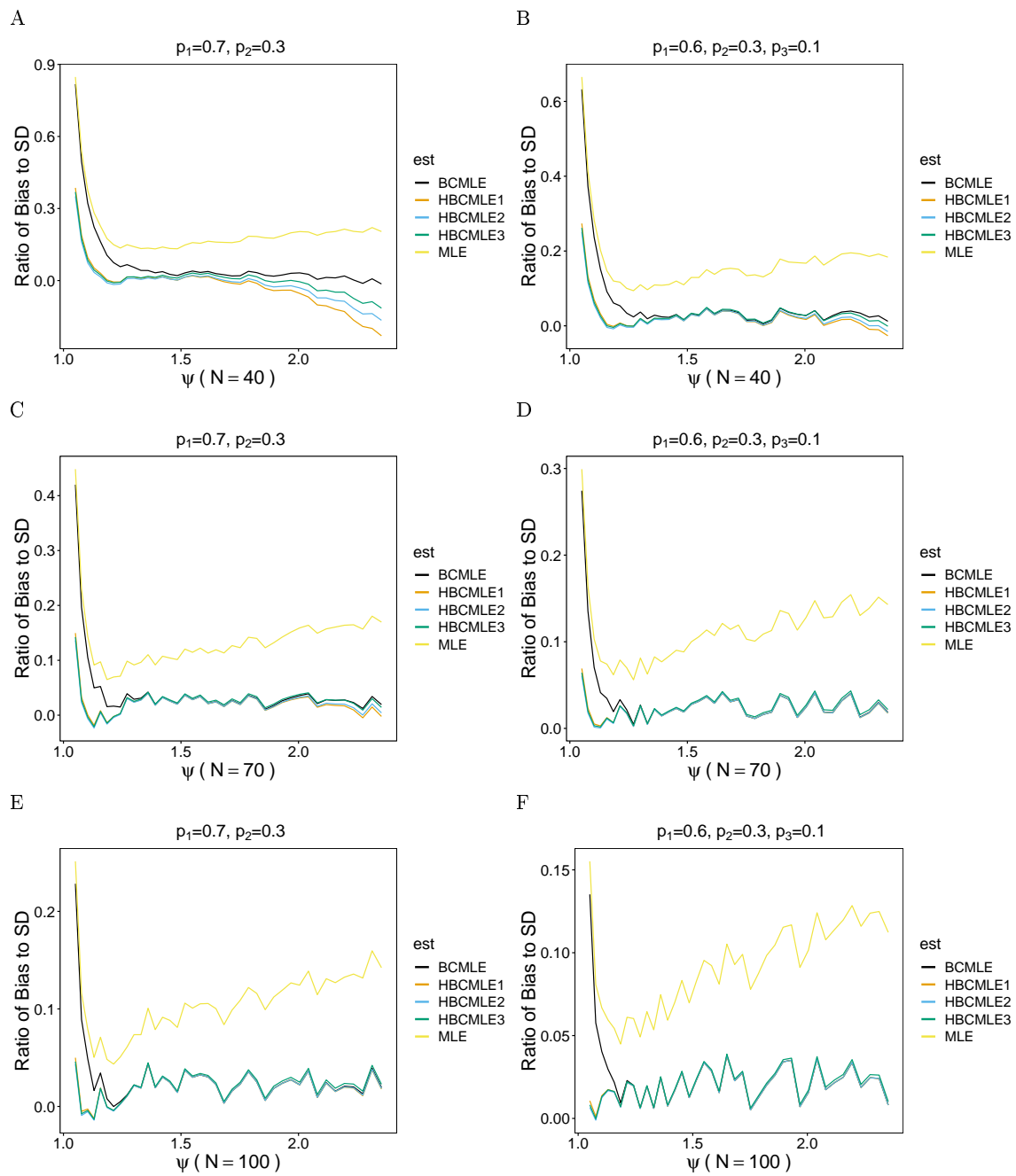


Figure 24: **Ratio of alternative estimators to BCMLE.** Similar to Figure 21

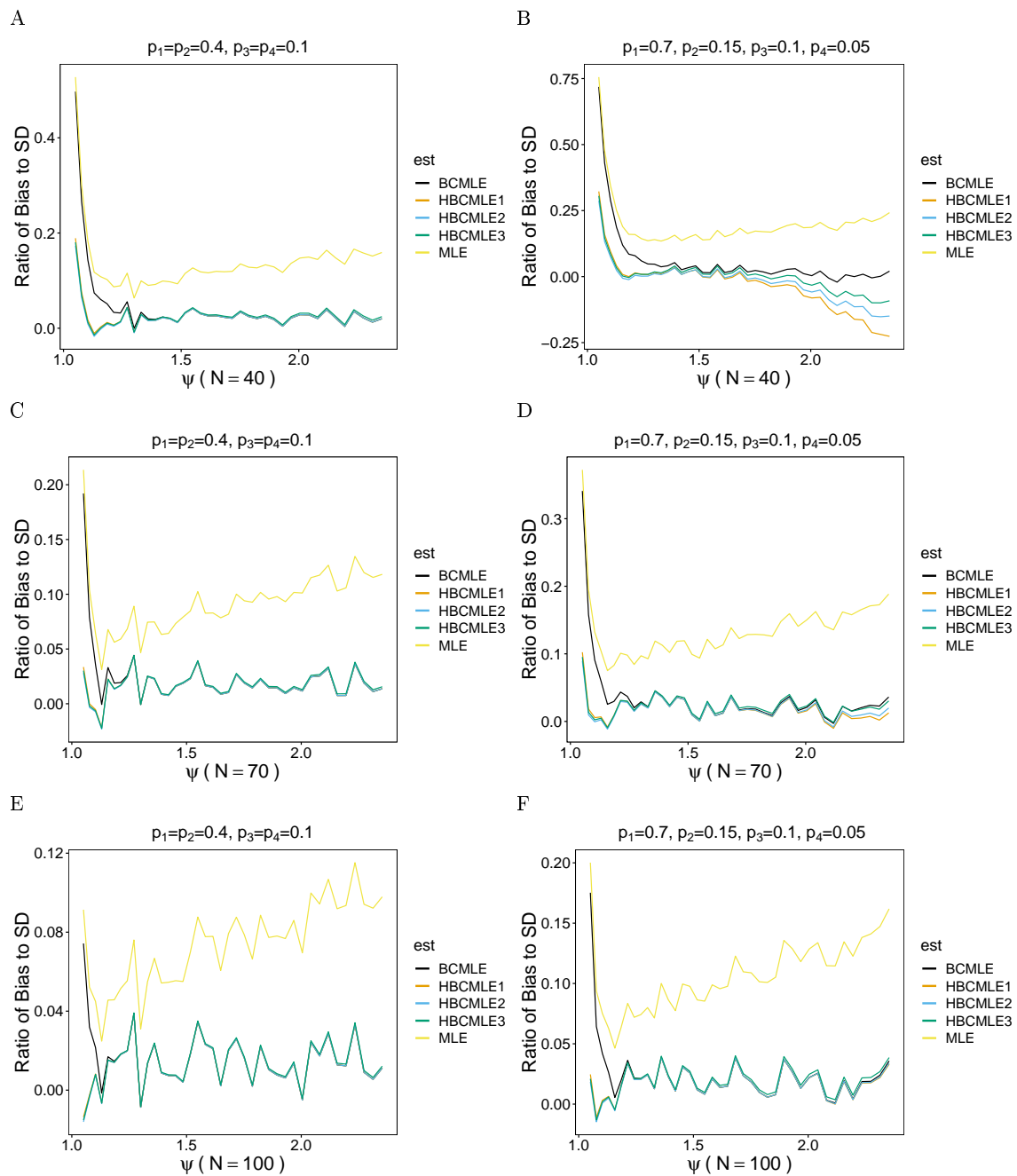


Figure 25: Same as Figure 21 but for different lineage-frequency distributions.

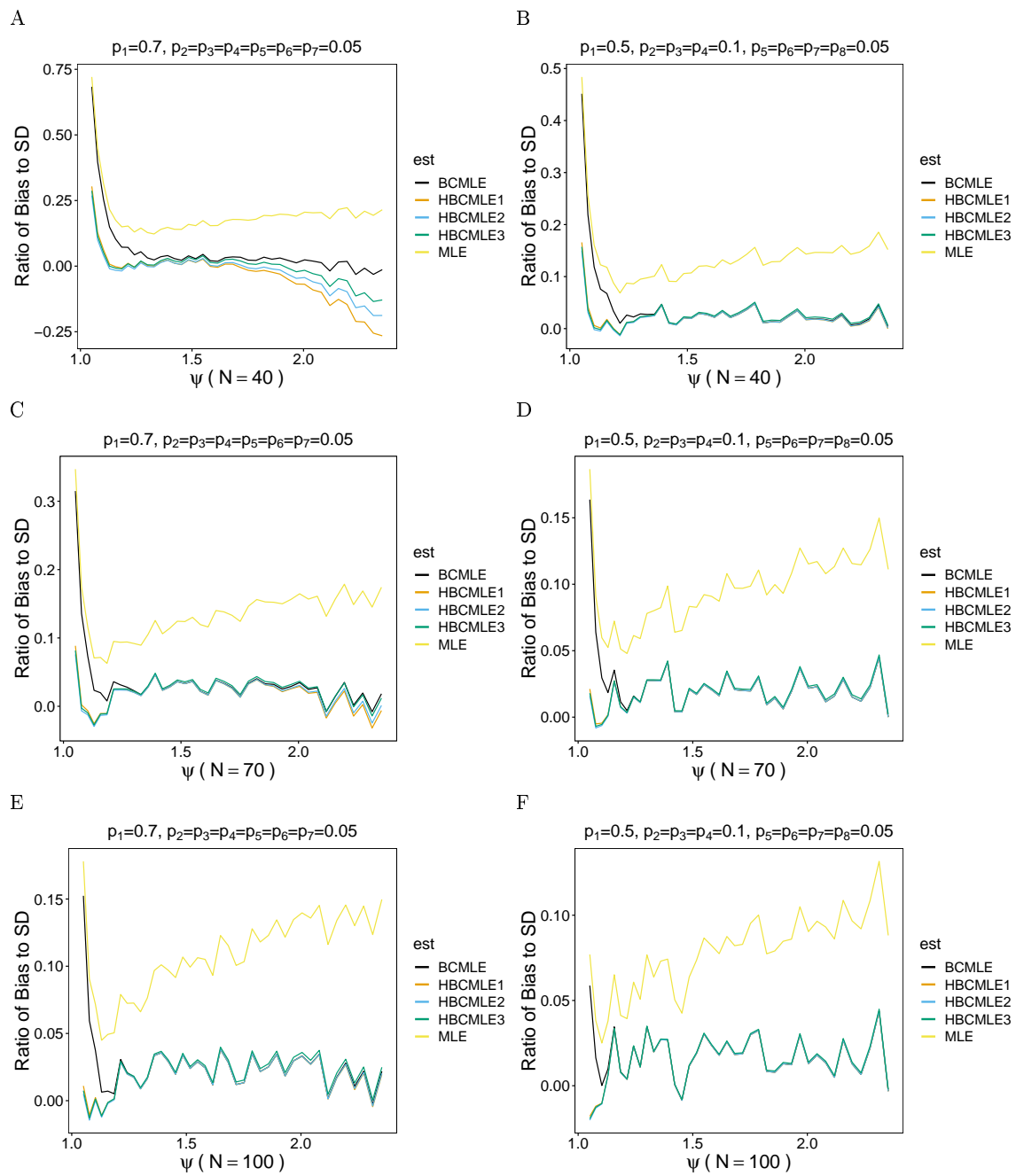


Figure 26: Same as Figure 21 but for different lineage-frequency distributions.



## 5 Histograms

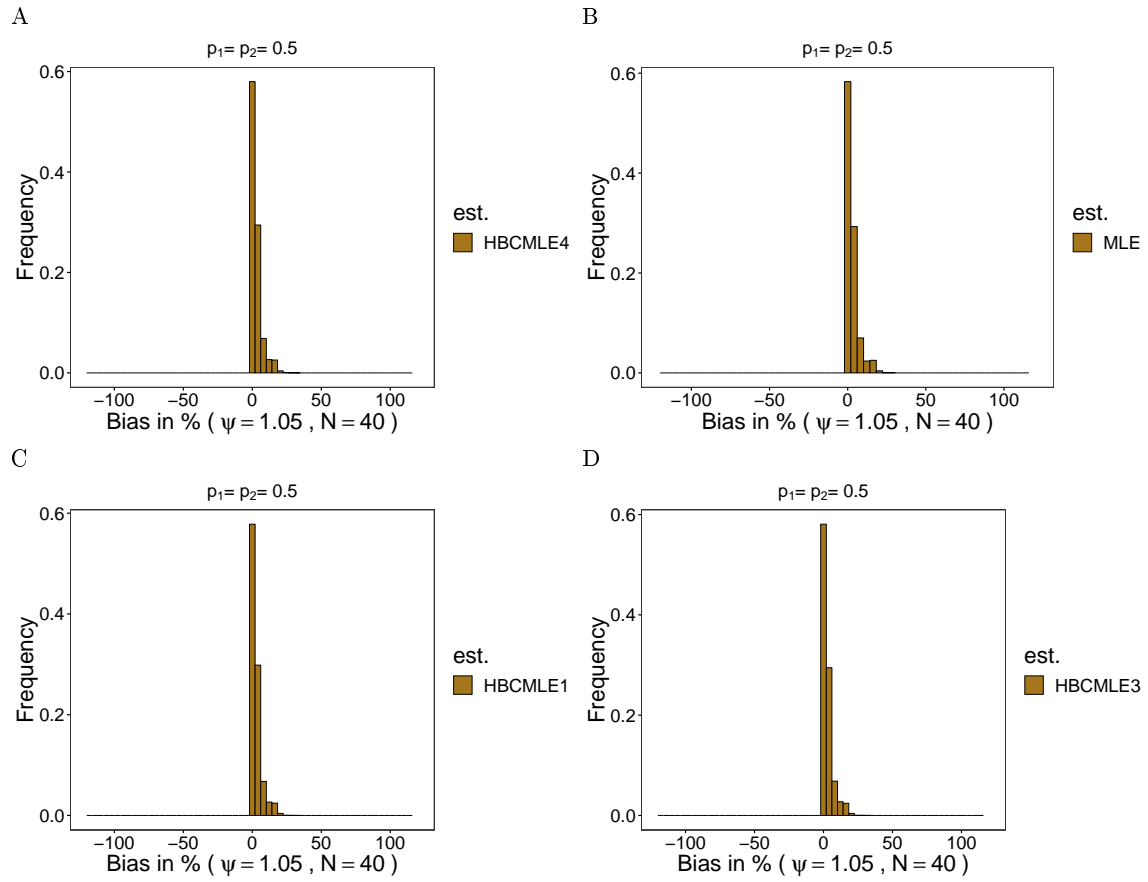


Figure 27: **Histograms.** The figure shows histograms of 10,000 estimations of  $\psi$  for a specific sample size  $N$  and true  $\psi$ . Each plot corresponds to a different estimator.

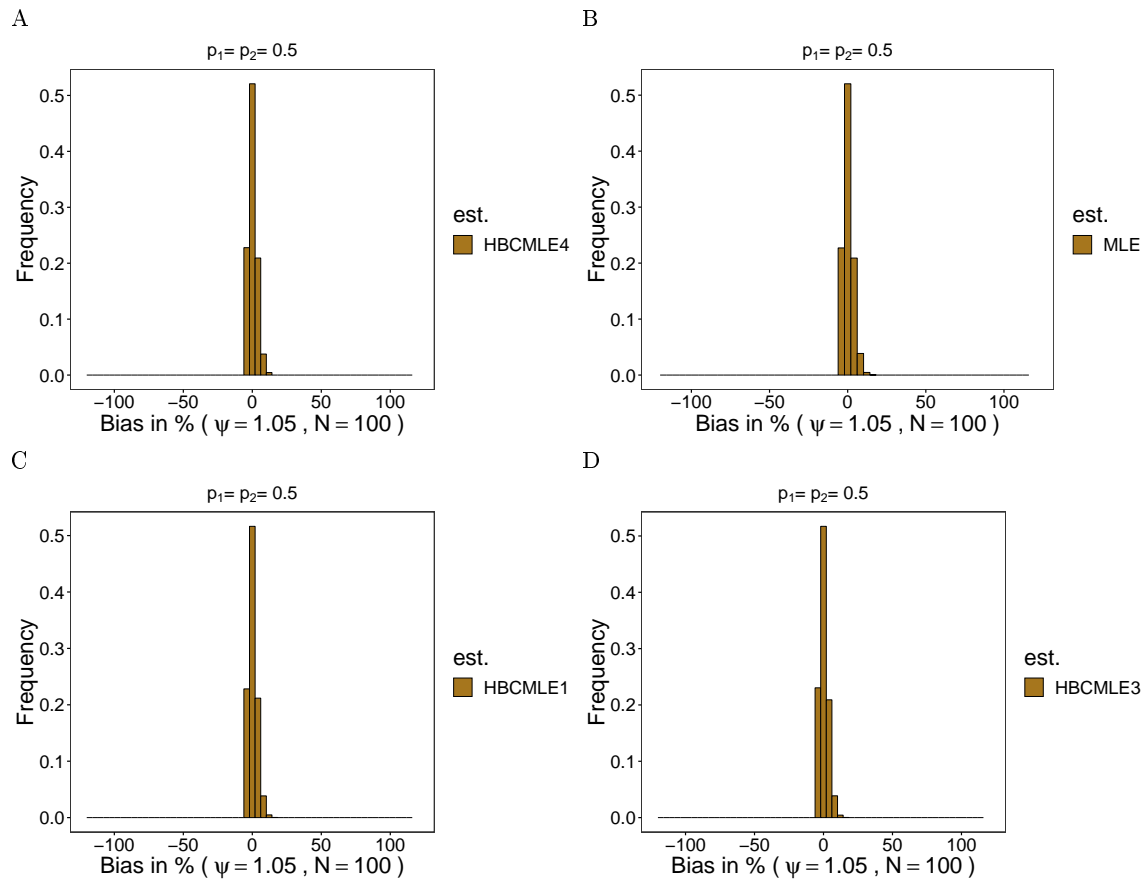


Figure 28: Same as Figure 27.

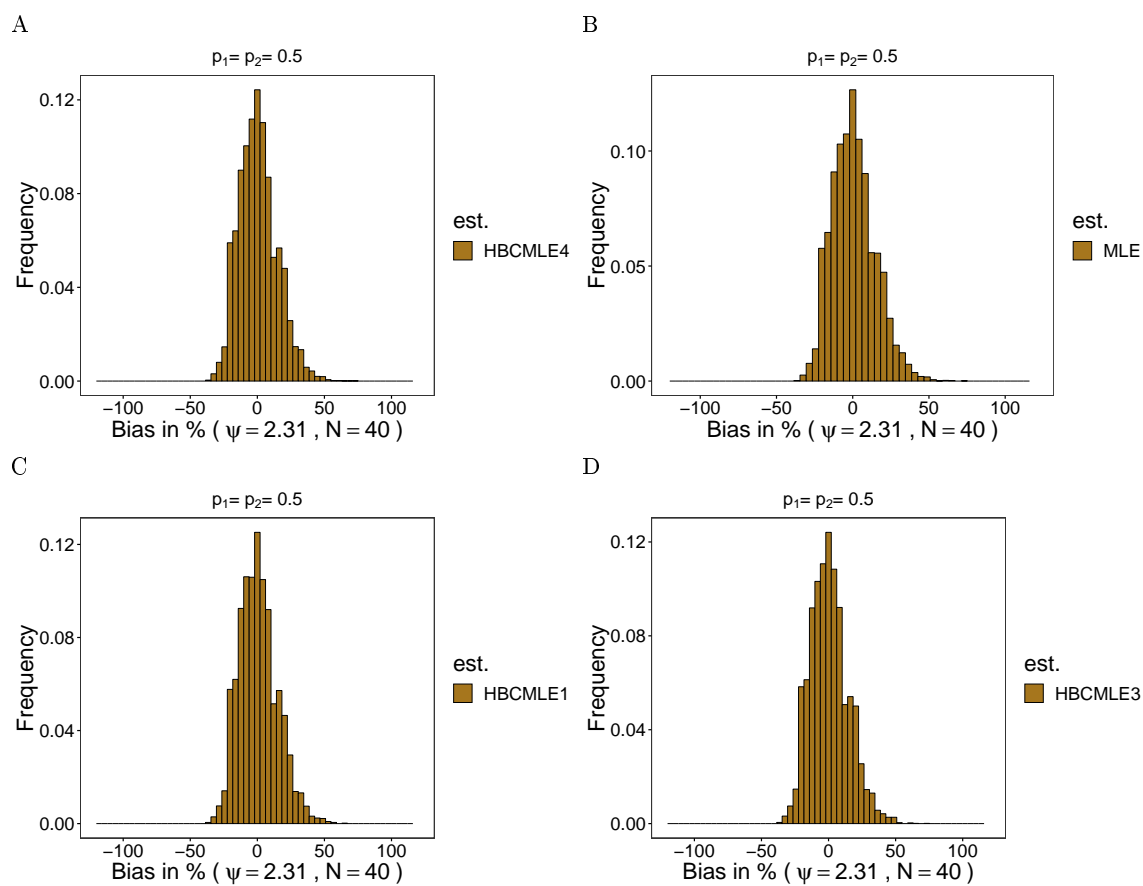


Figure 29: Same as Figure 27.

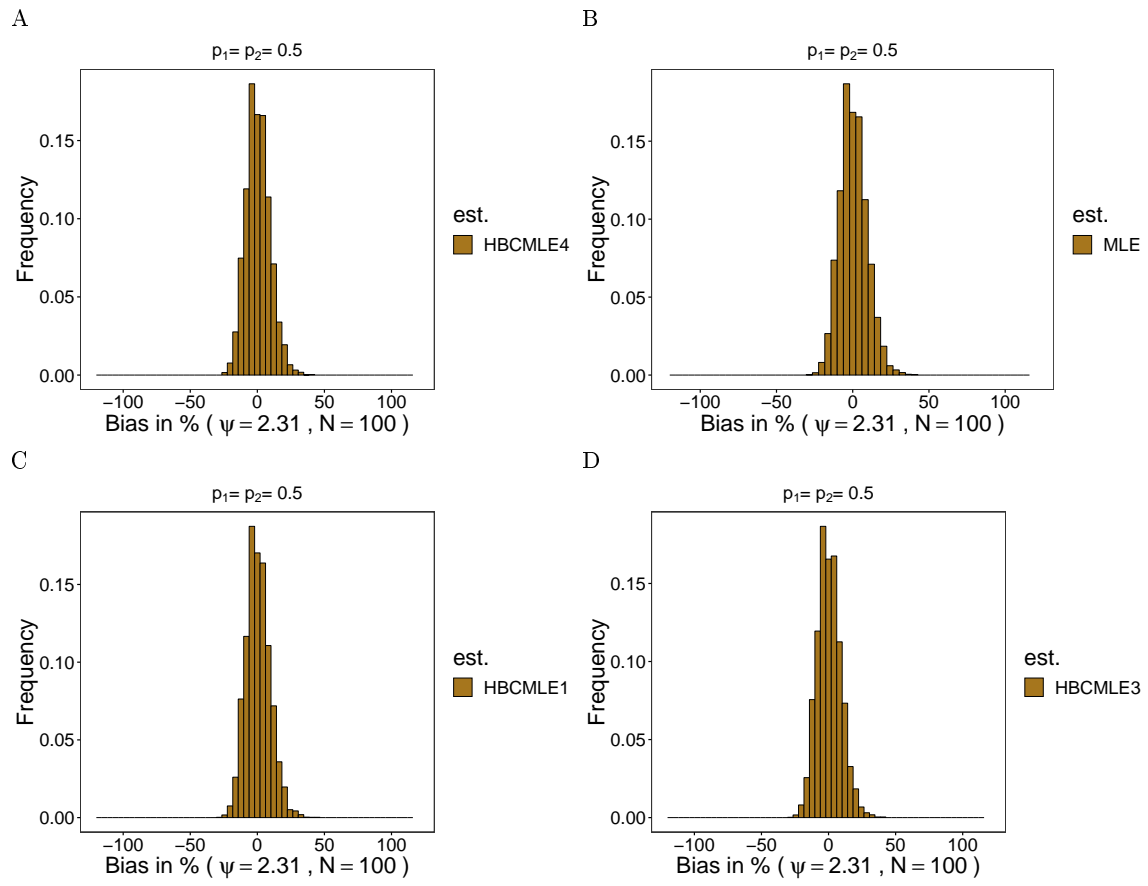


Figure 30: Same as Figure 27.

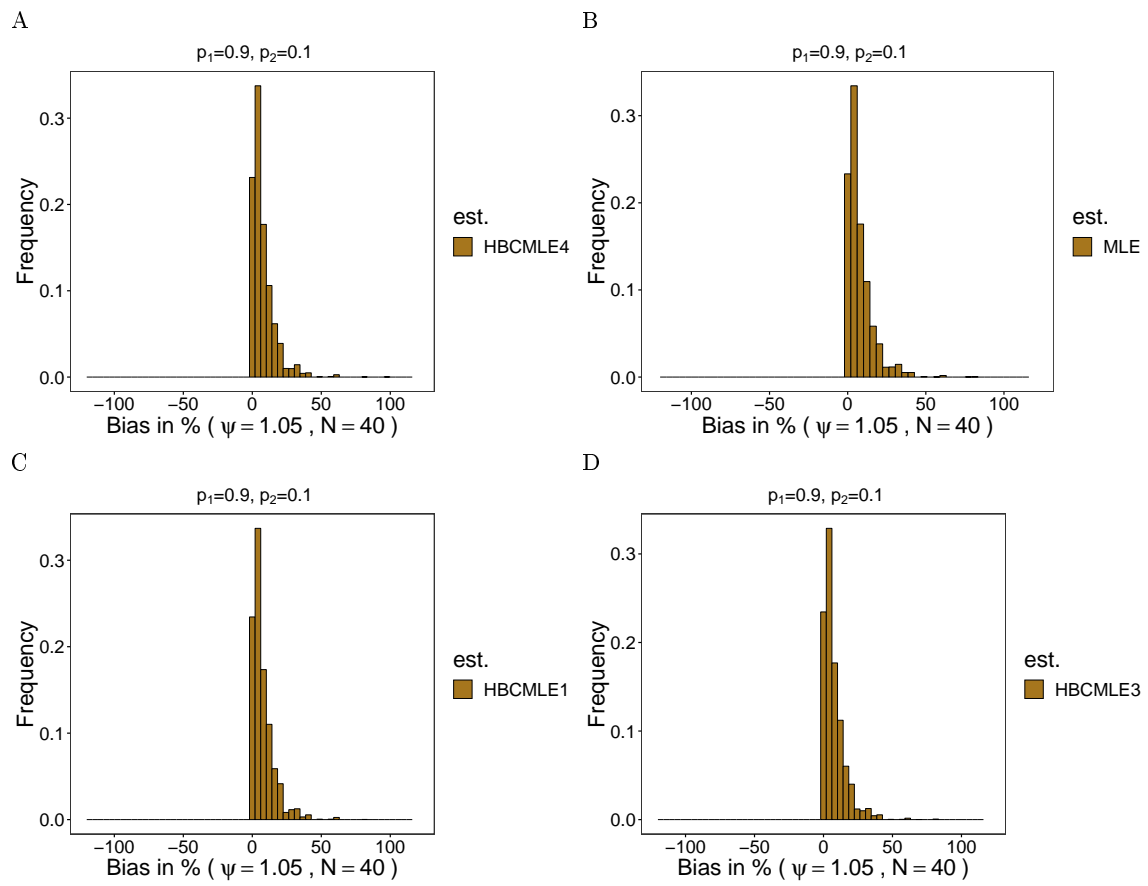


Figure 31: Same as Figure 27.

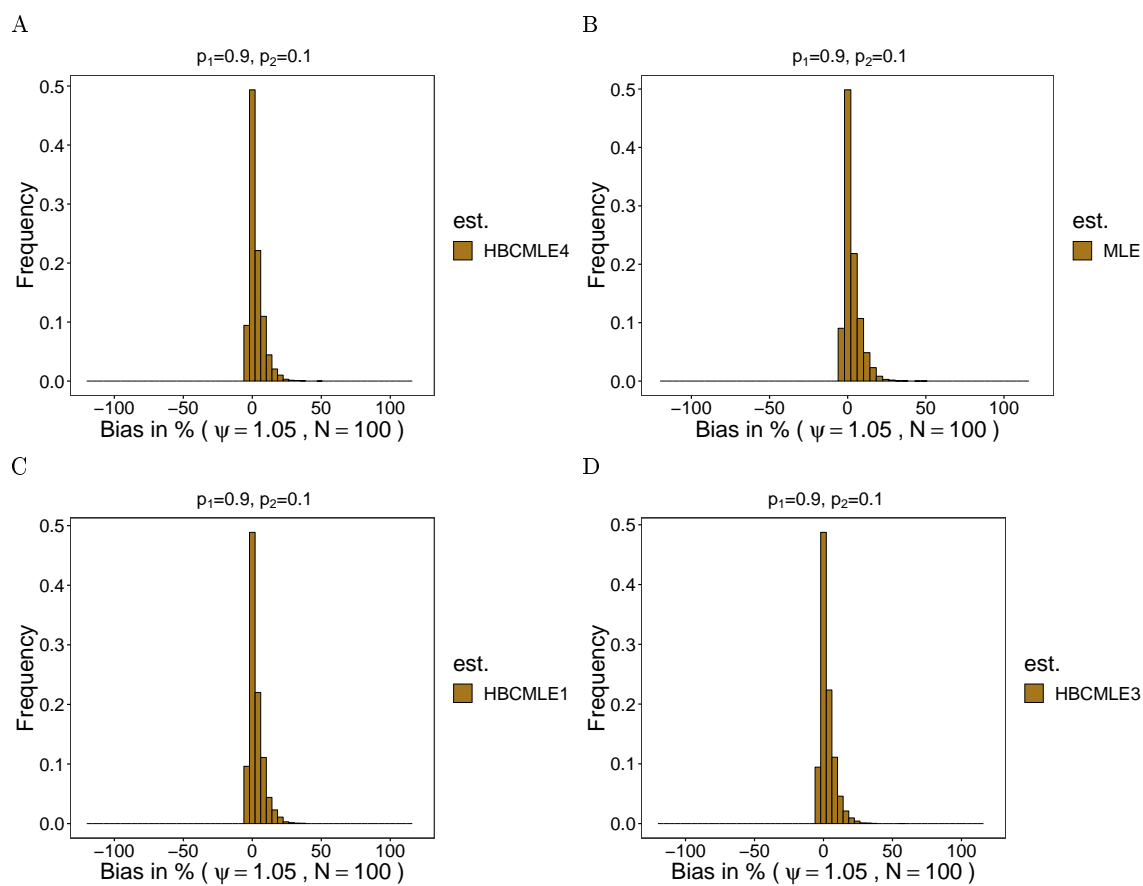


Figure 32: Same as Figure 27.

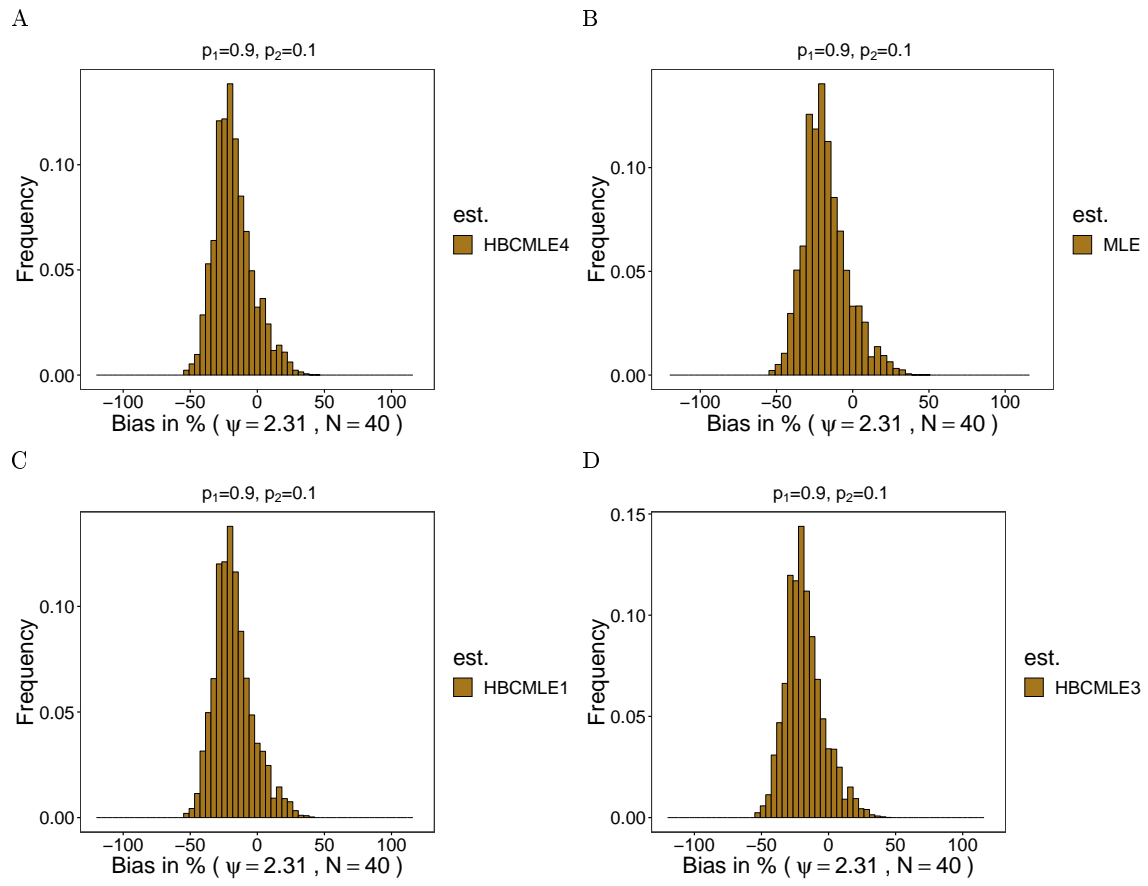


Figure 33: Same as Figure 27.

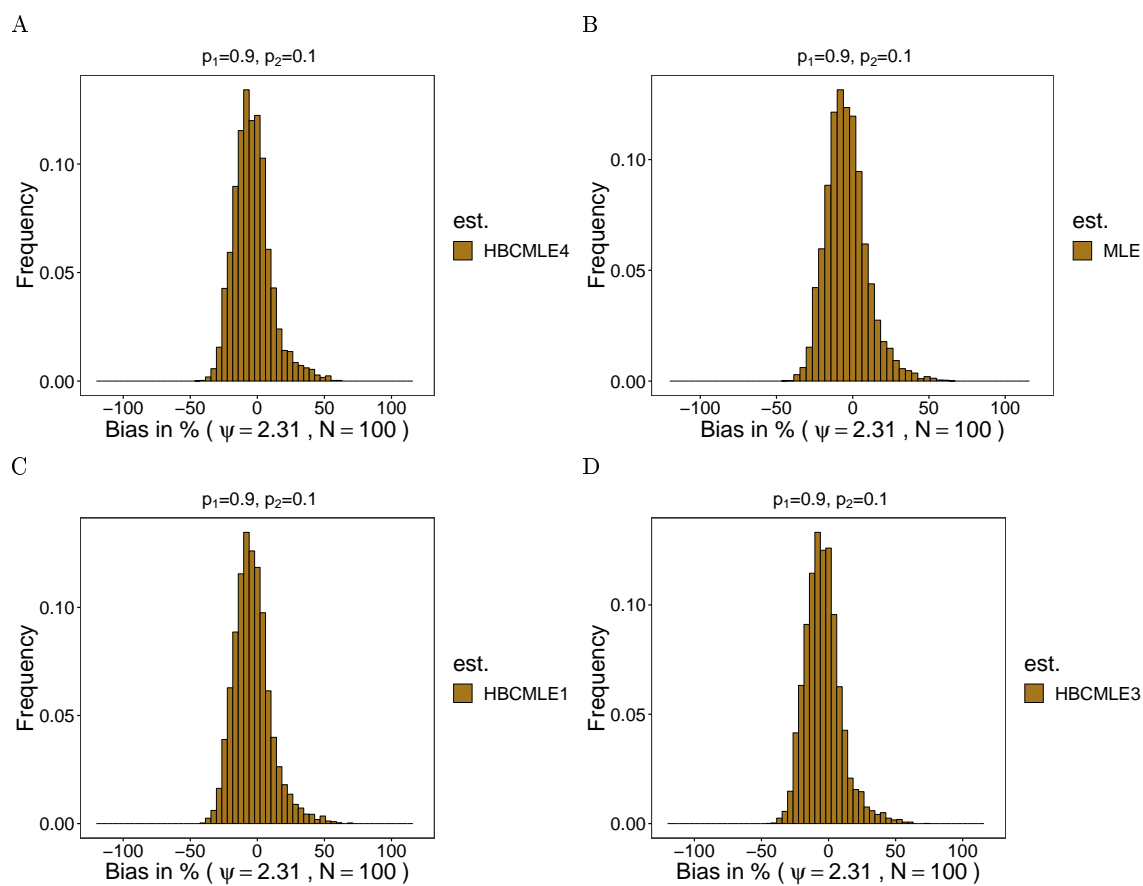


Figure 34: Same as Figure 27.



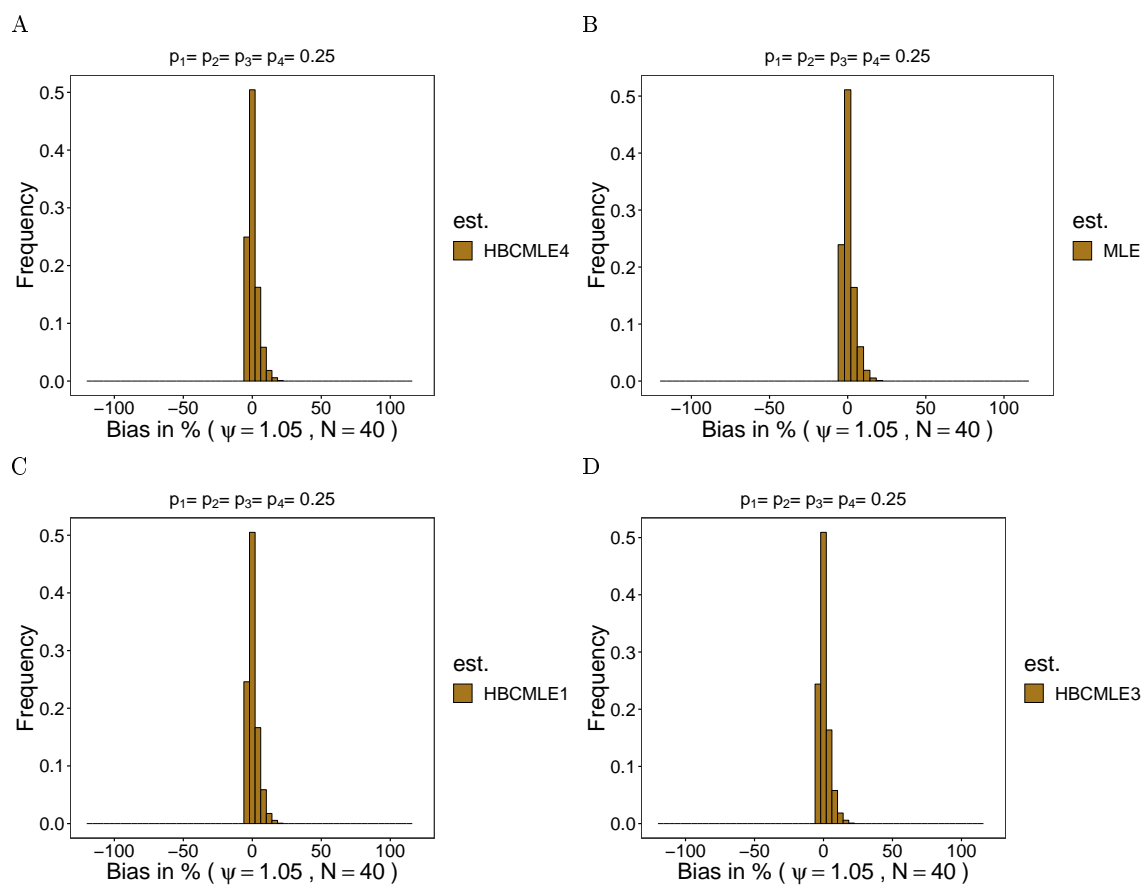


Figure 35: Same as Figure 27.

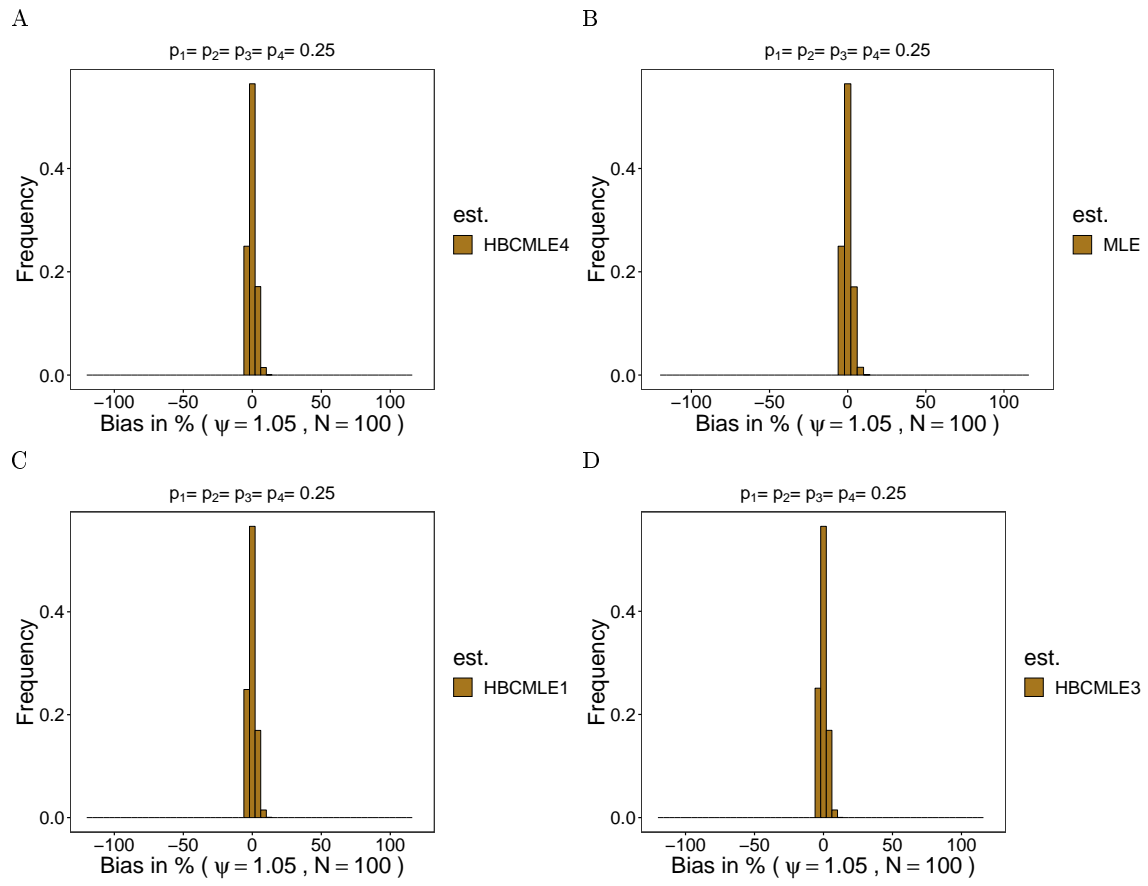


Figure 36: Same as Figure 27.

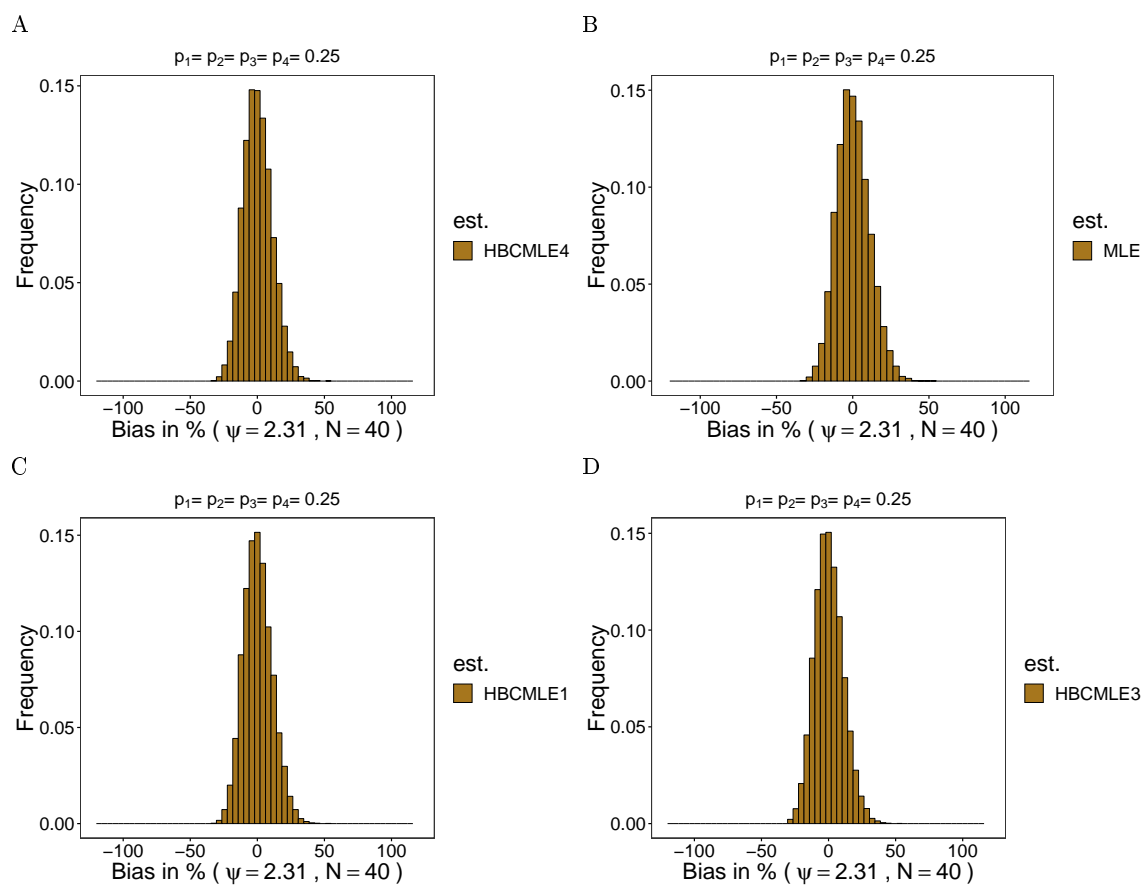


Figure 37: Same as Figure 27.

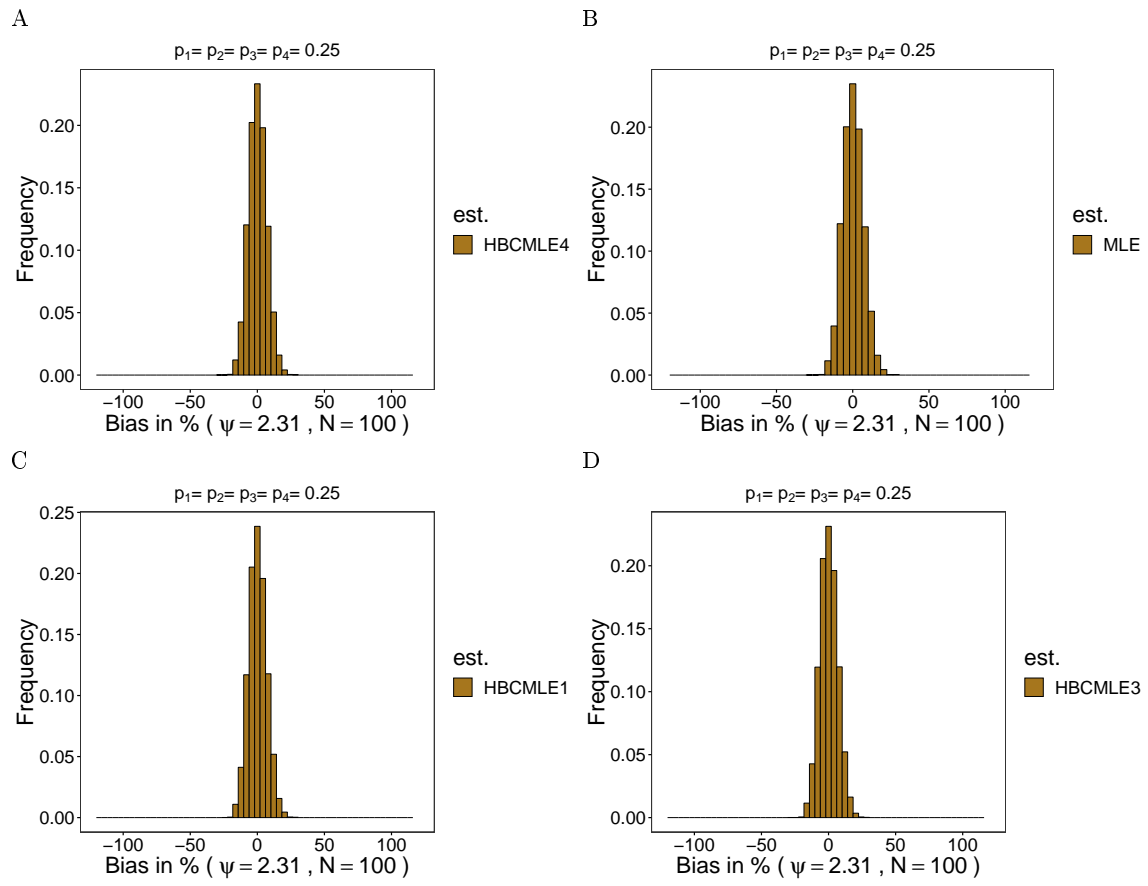


Figure 38: Same as Figure 27.

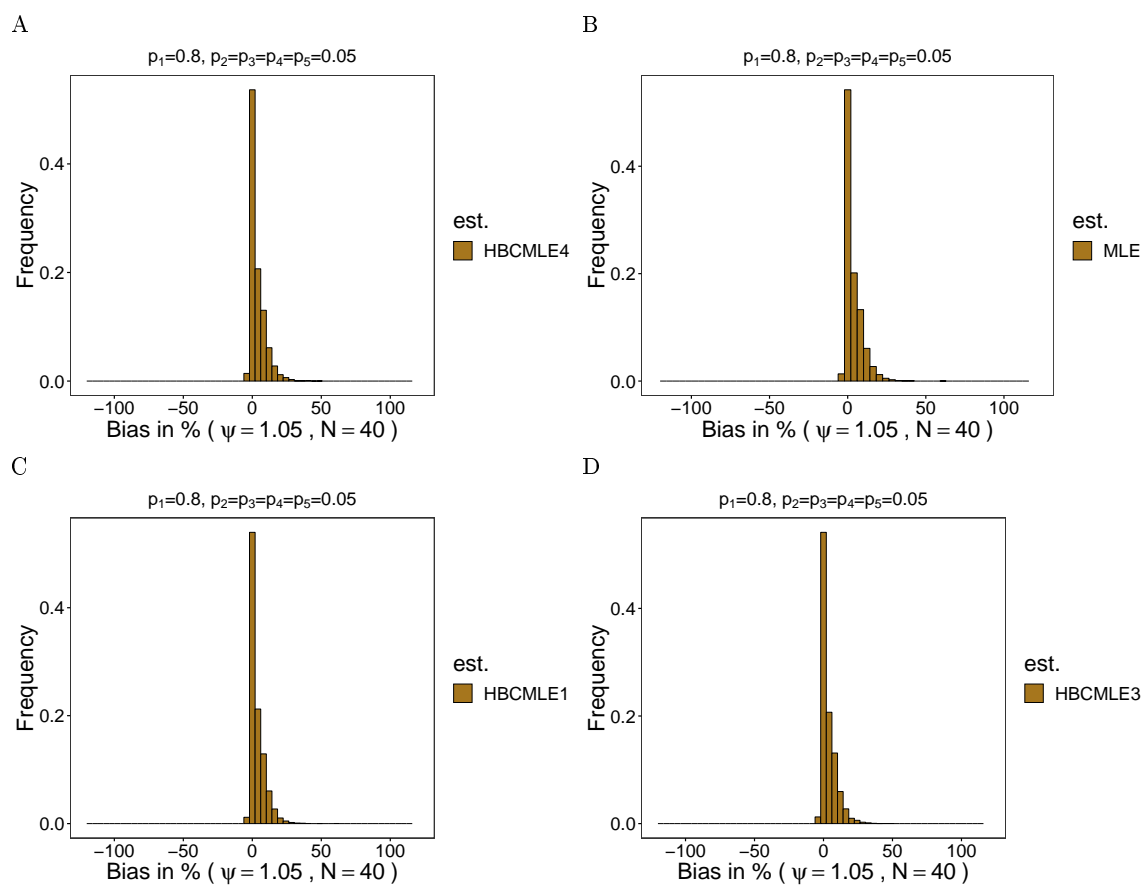


Figure 39: Same as Figure 27.

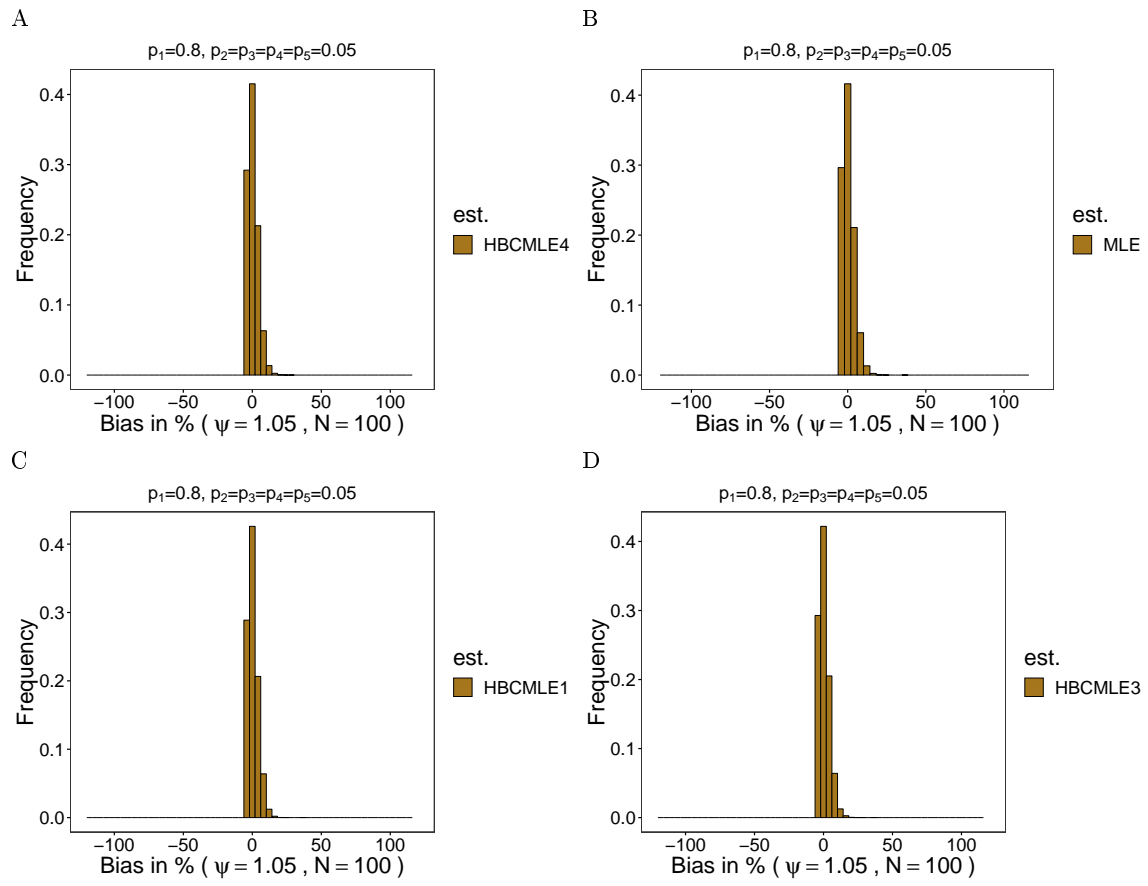


Figure 40: Same as Figure 27.

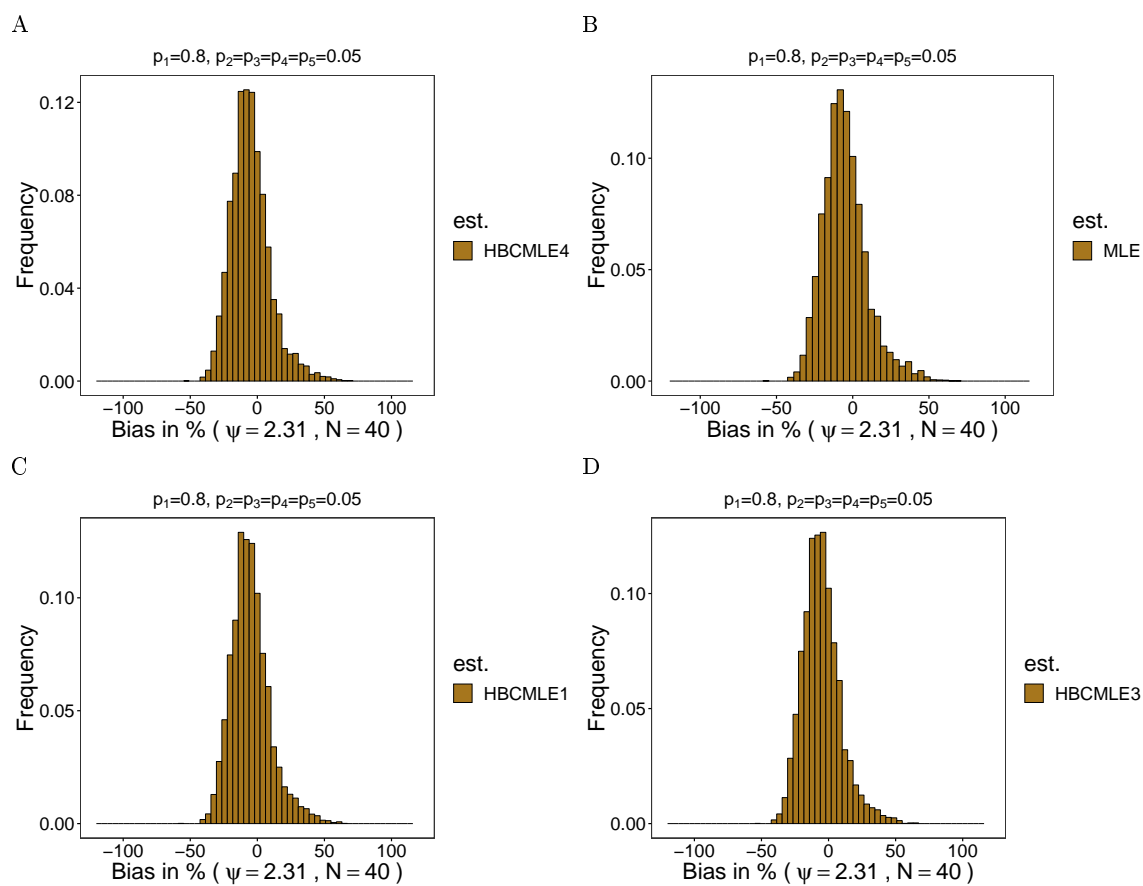


Figure 41: Same as Figure 27.

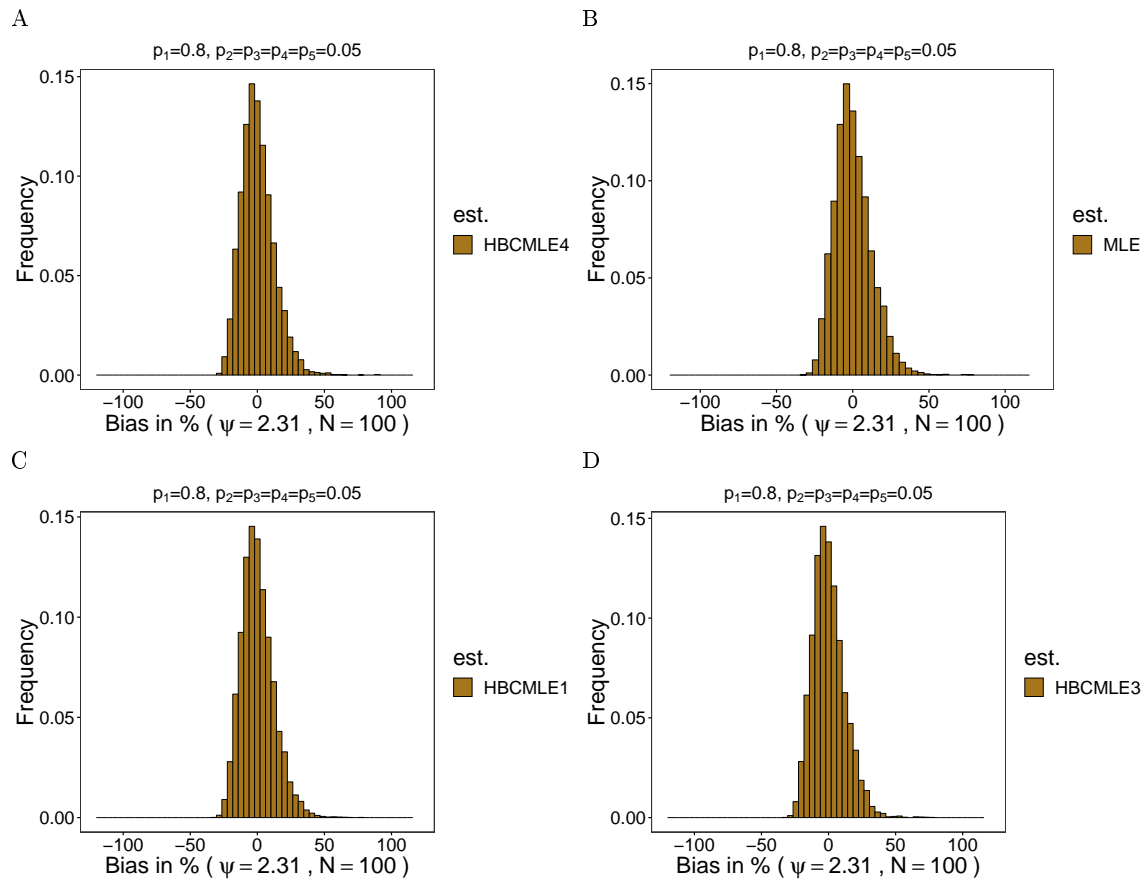


Figure 42: Same as Figure 27.



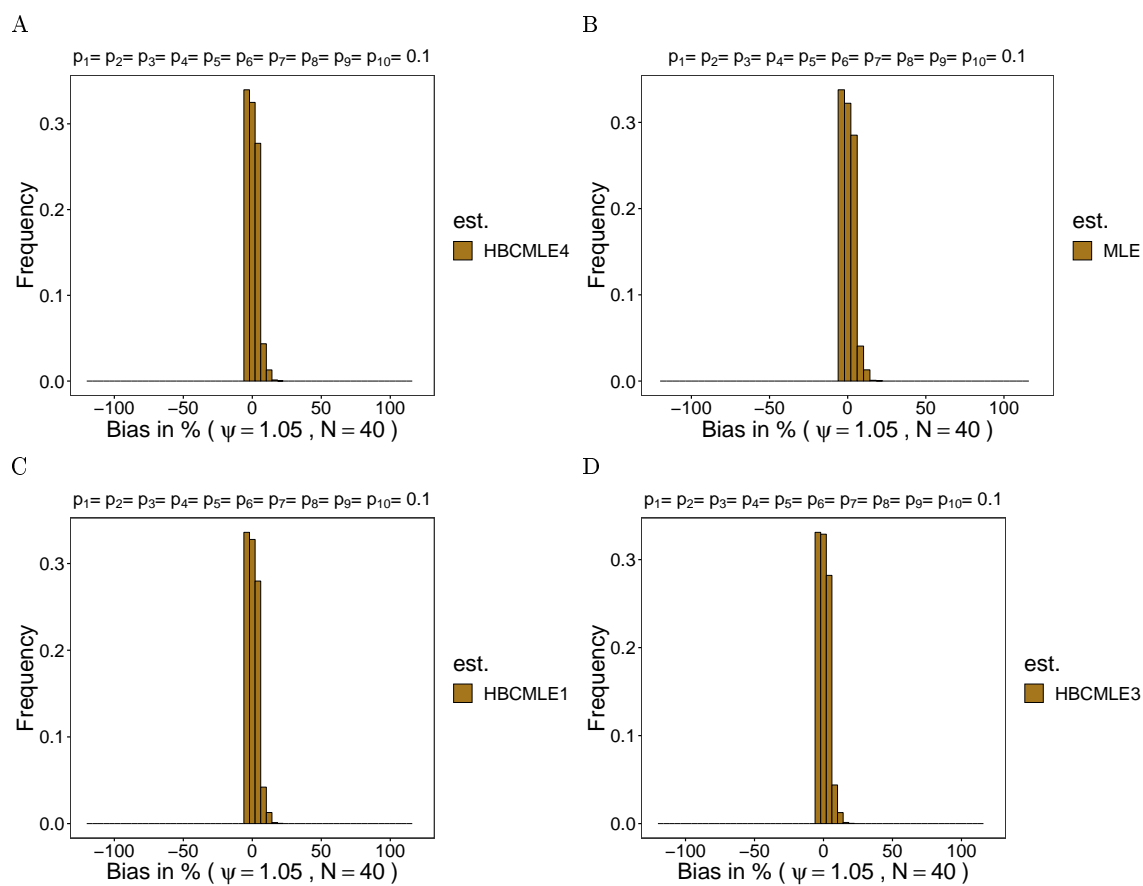


Figure 43: Same as Figure 27.

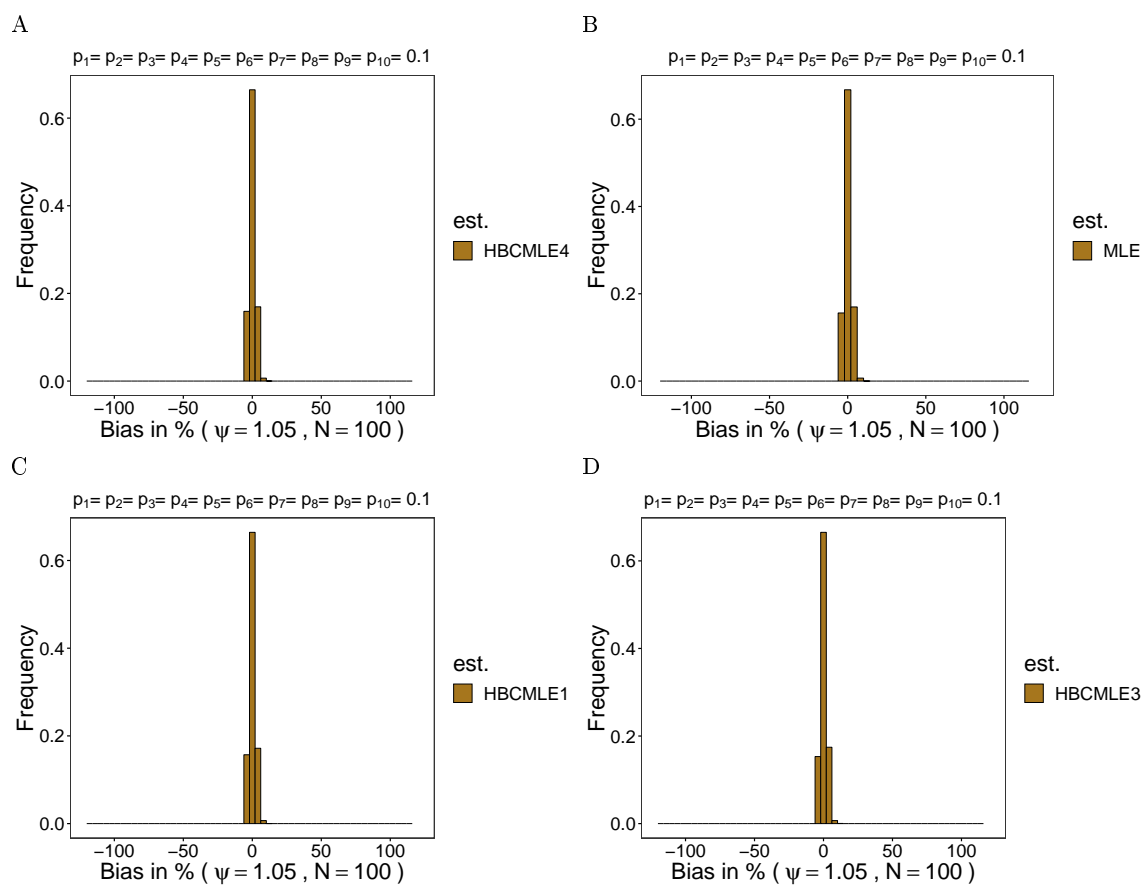


Figure 44: Same as Figure 27.

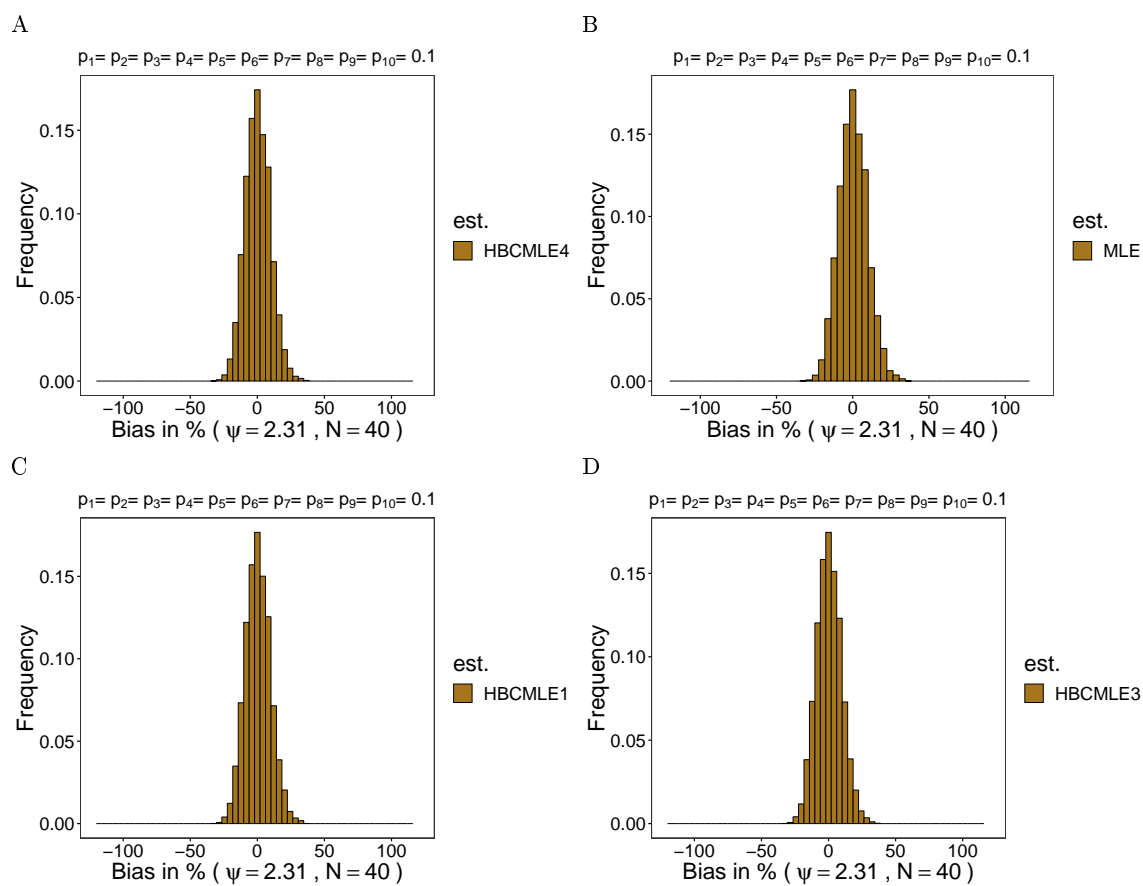


Figure 45: Same as Figure 27.

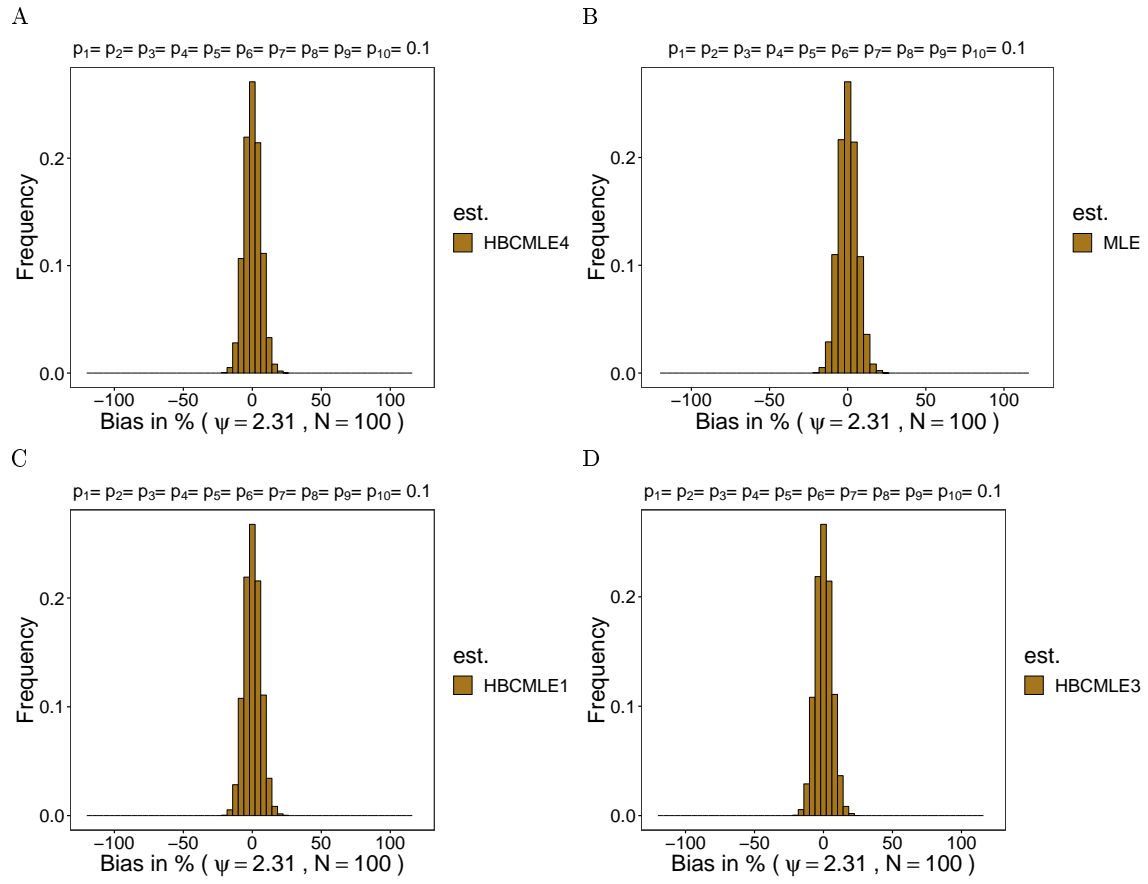


Figure 46: Same as Figure 27.

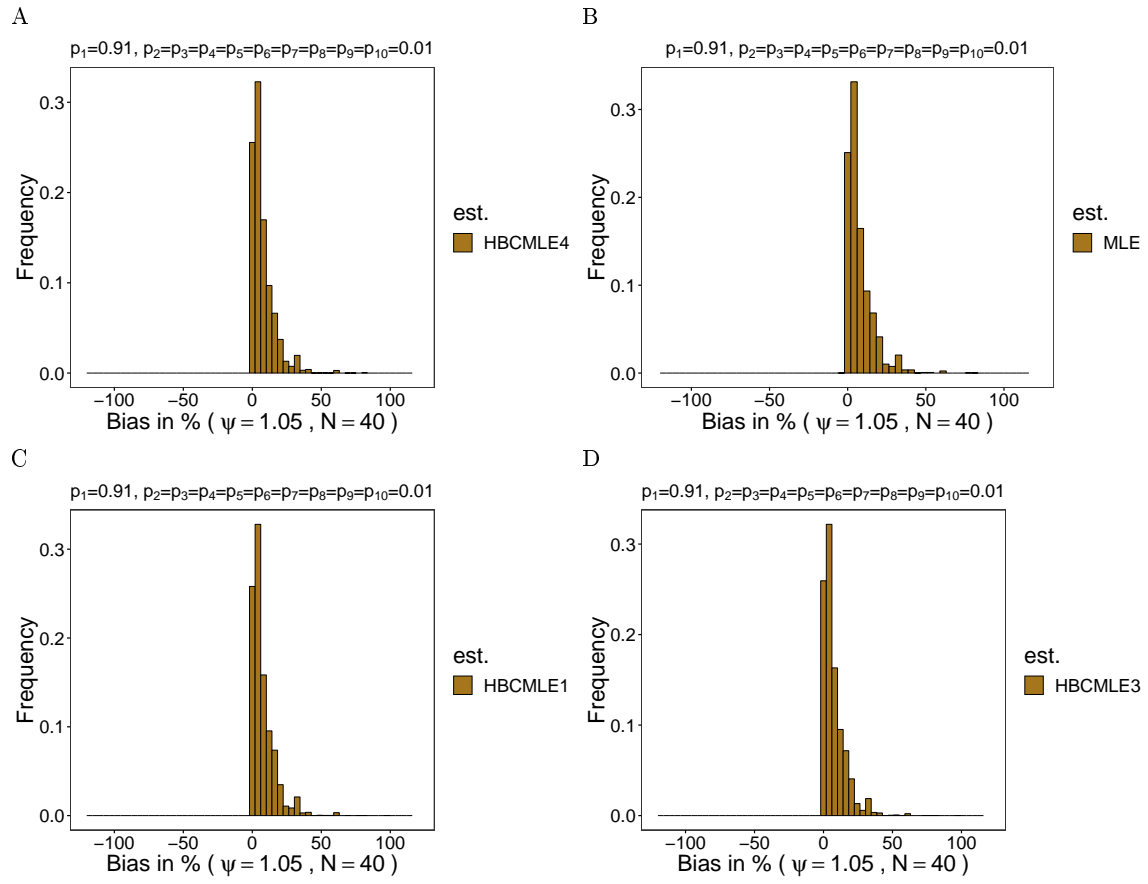


Figure 47: Same as Figure 27.

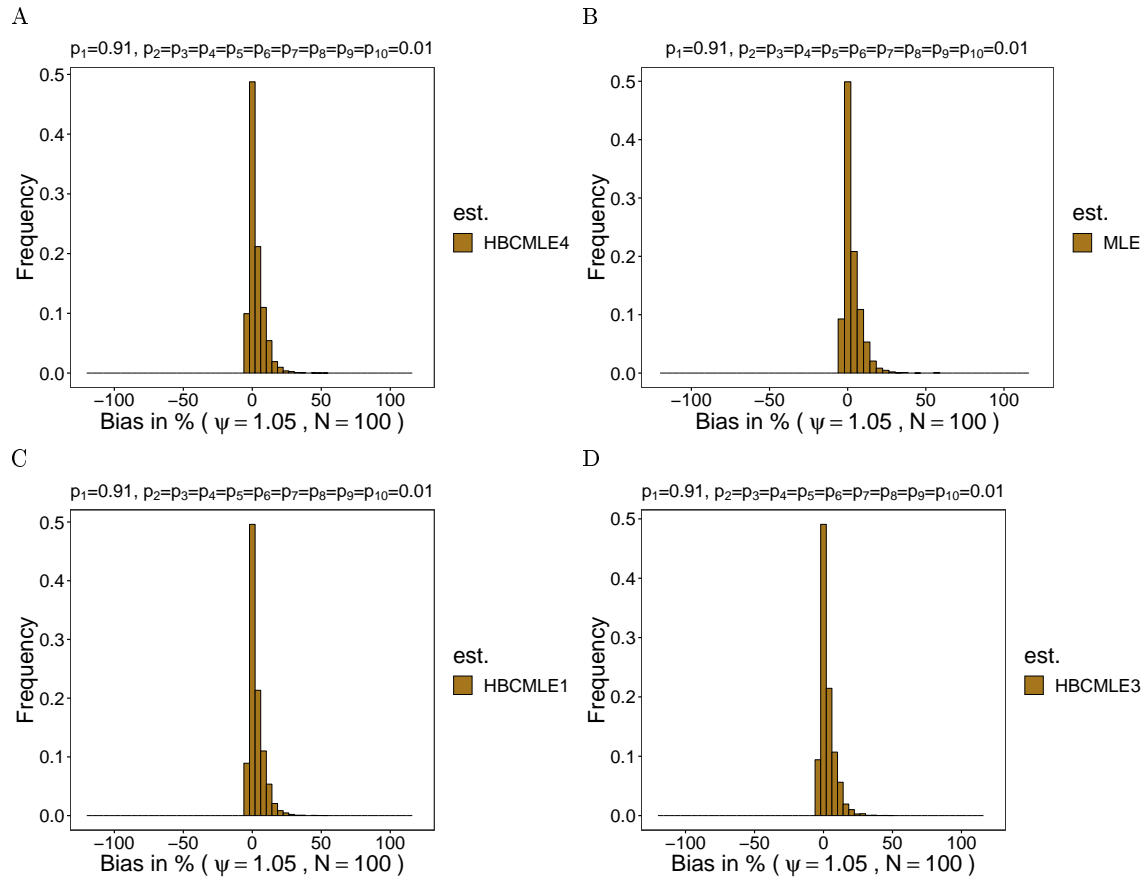


Figure 48: Same as Figure 27.

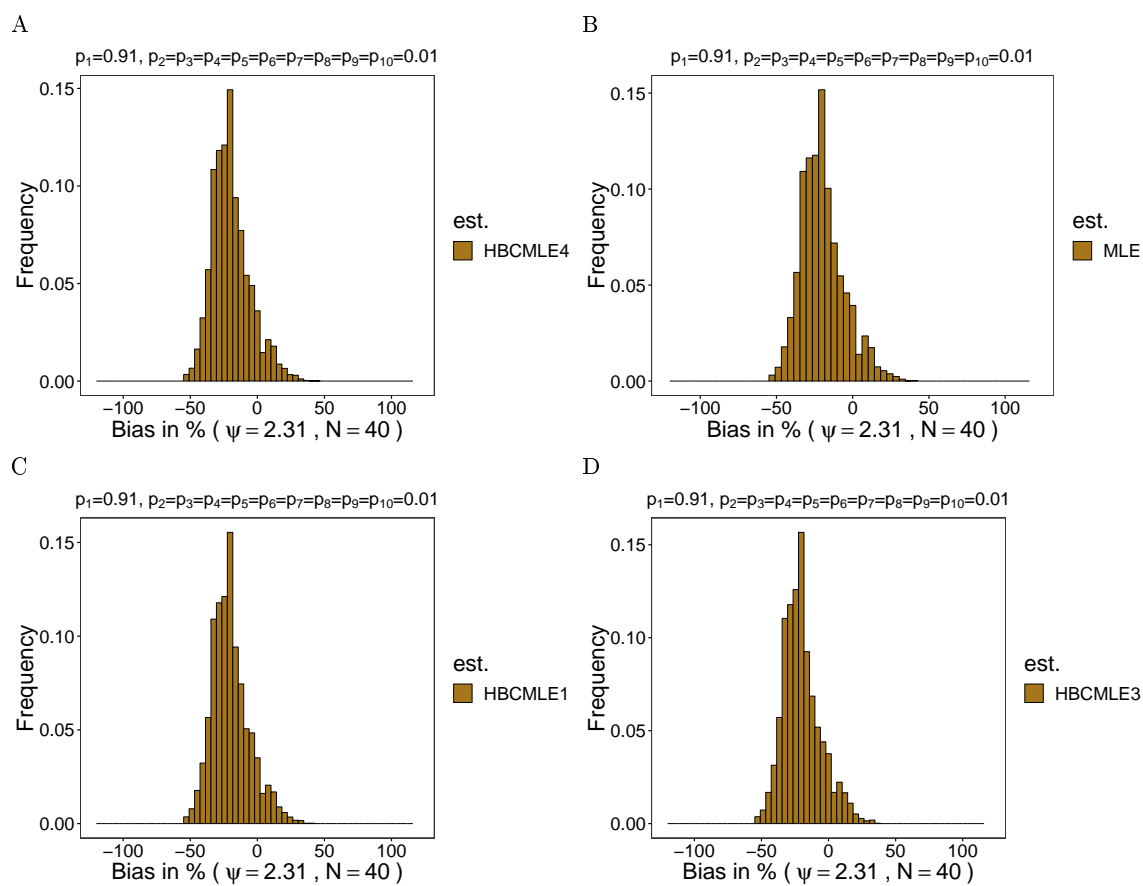


Figure 49: Same as Figure 27.

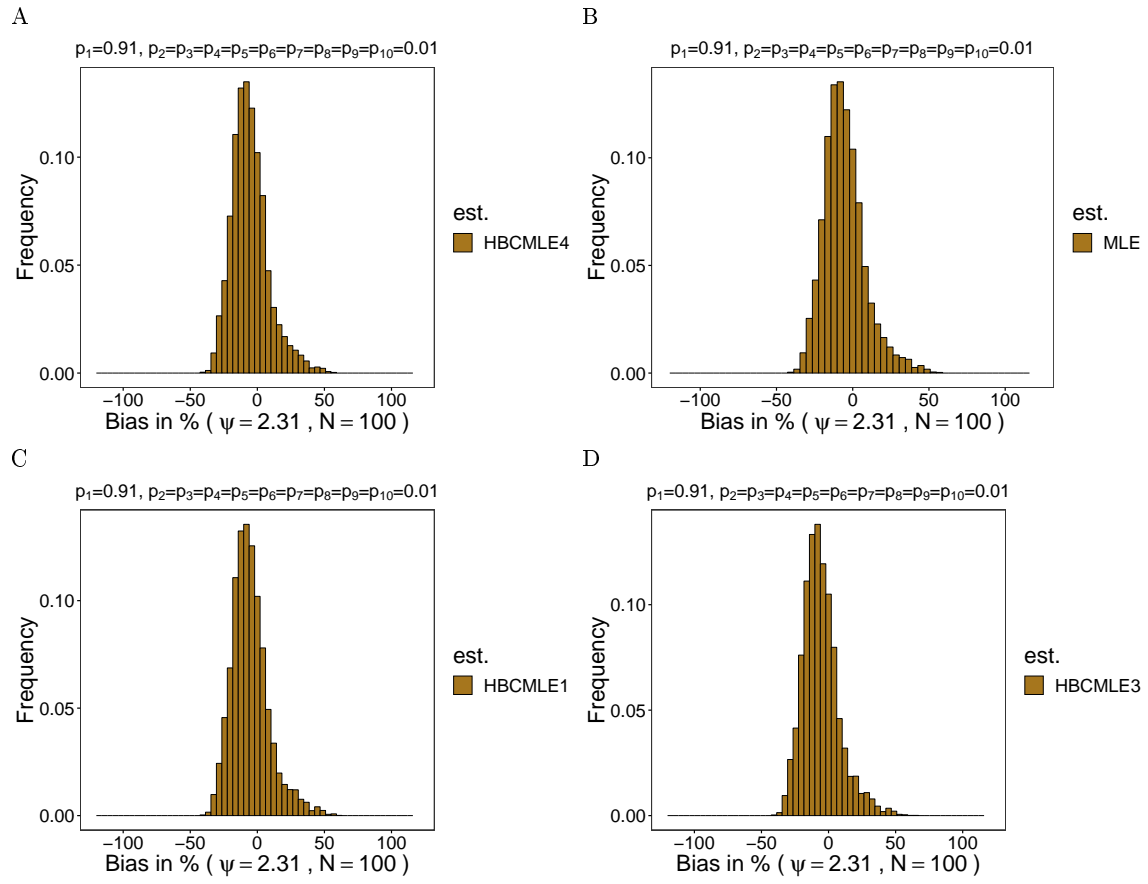


Figure 50: Same as Figure 27.





## 6 Boxplots

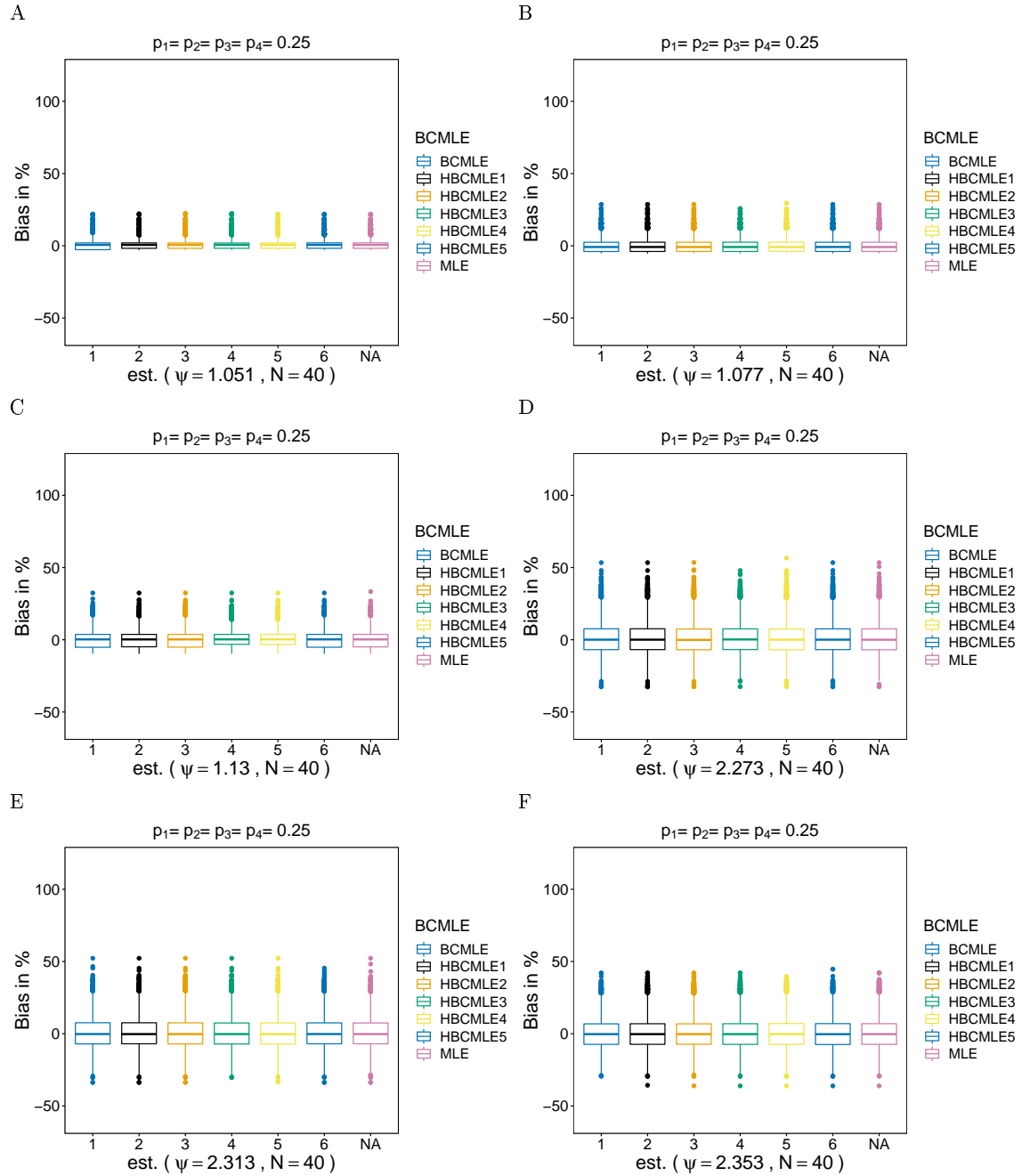


Figure 51: **Boxplots.** The figure shows boxplots of 10,000 estimations of  $\psi$  for a specific sample size  $N$ . Colors correspond to different estimators. Each plot corresponds to a different value of true parameter  $\psi$ .

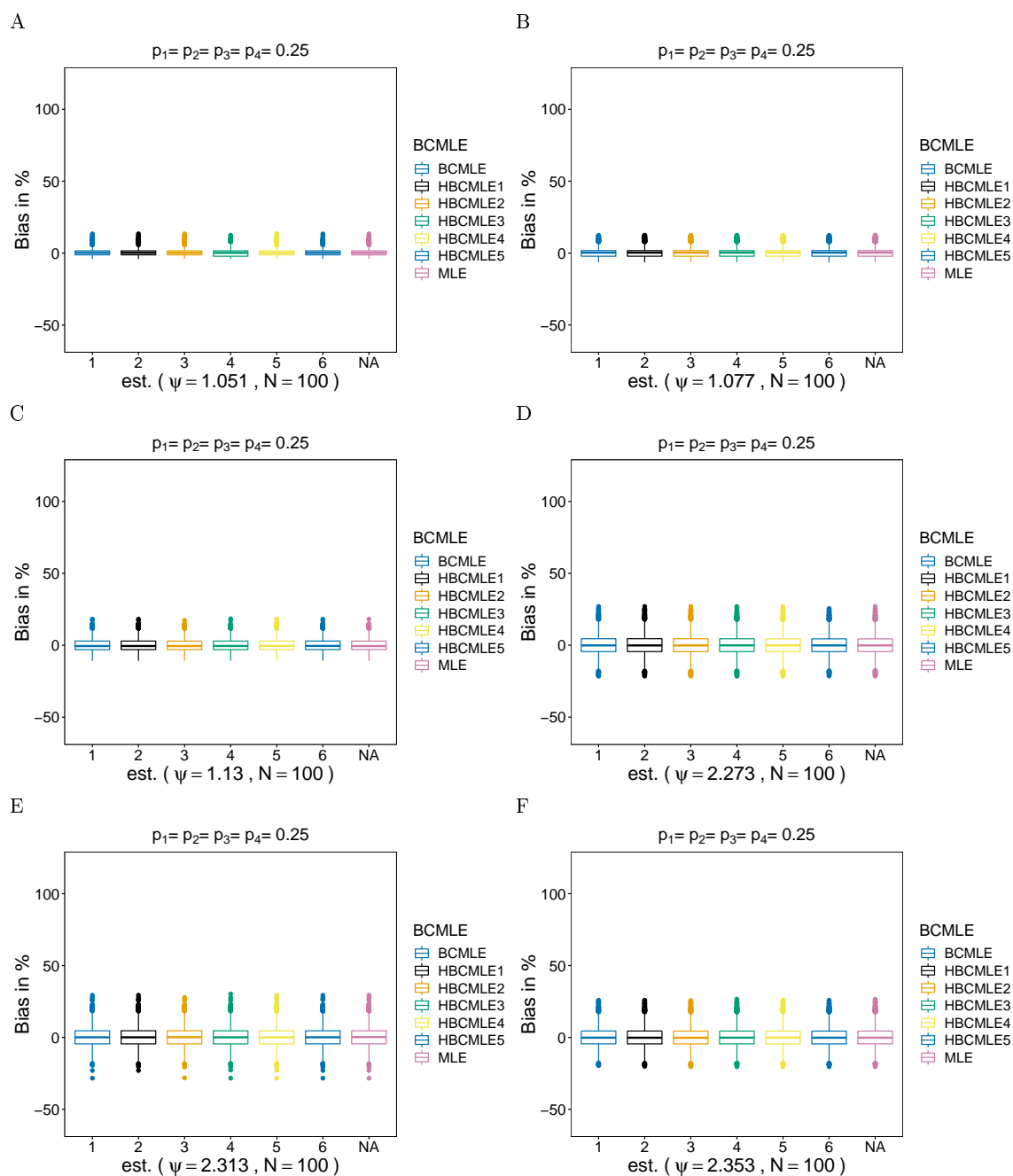


Figure 52: Same as Figure 51

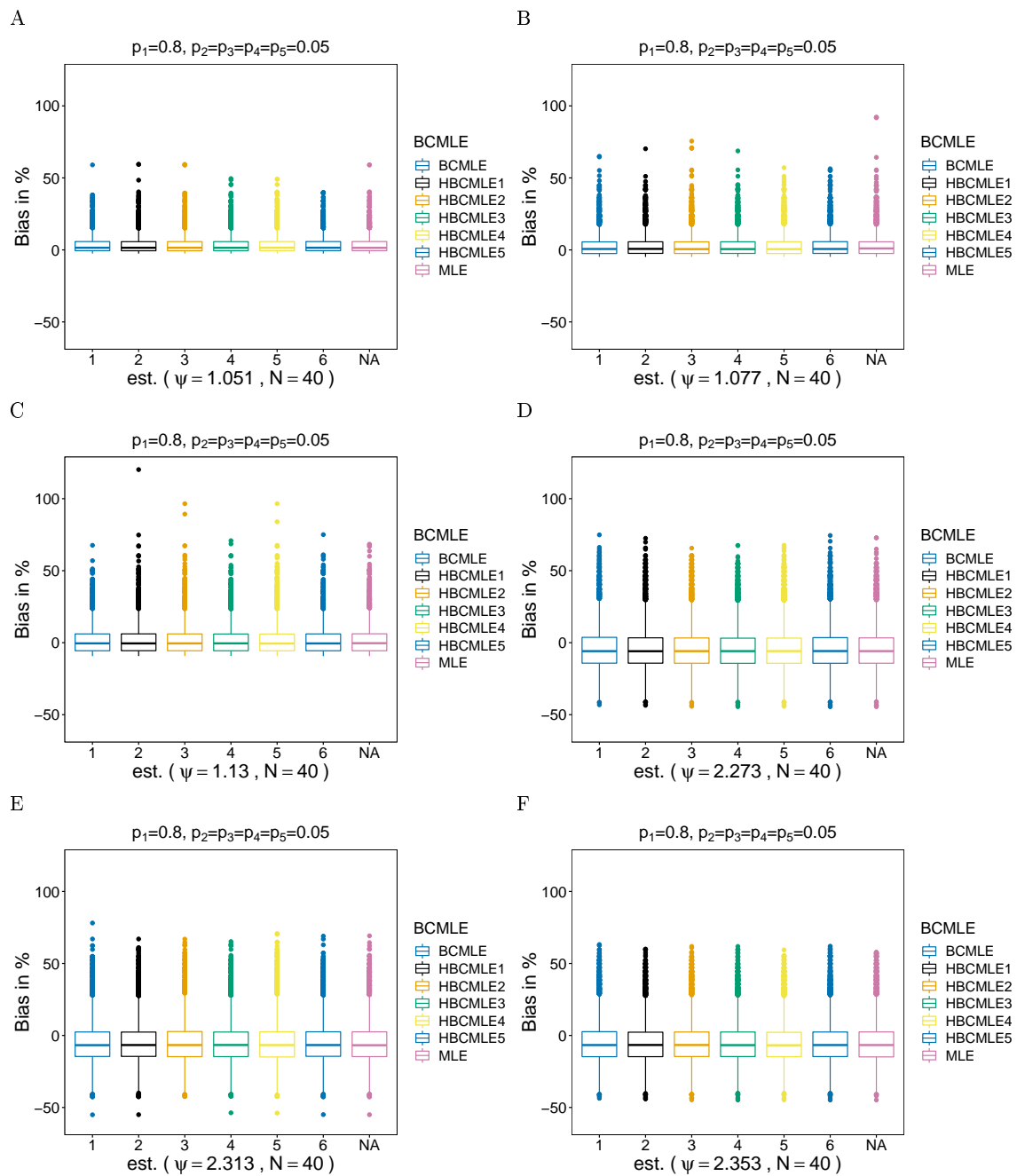


Figure 53: Same as Figure 51

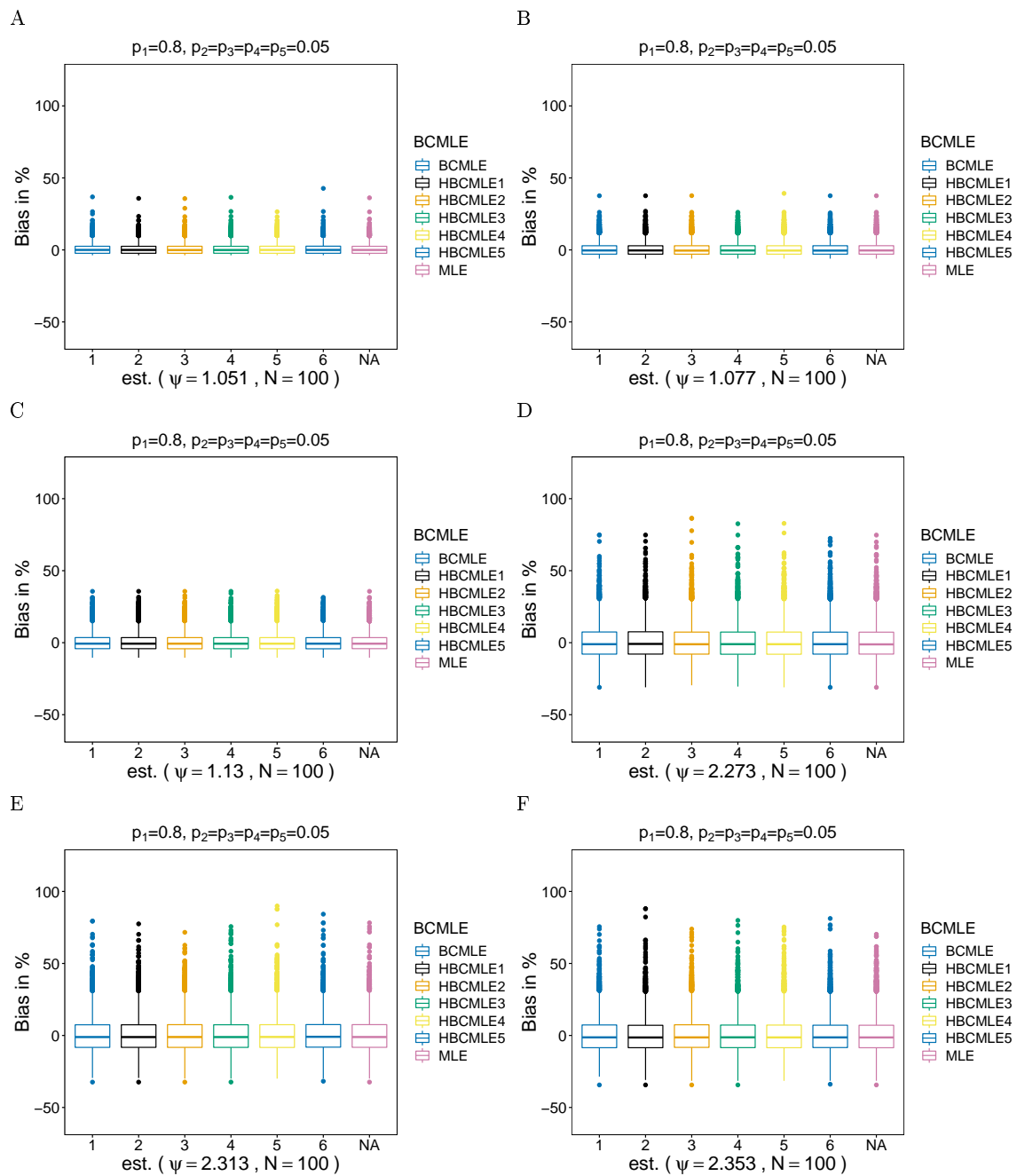


Figure 54: Same as Figure 51

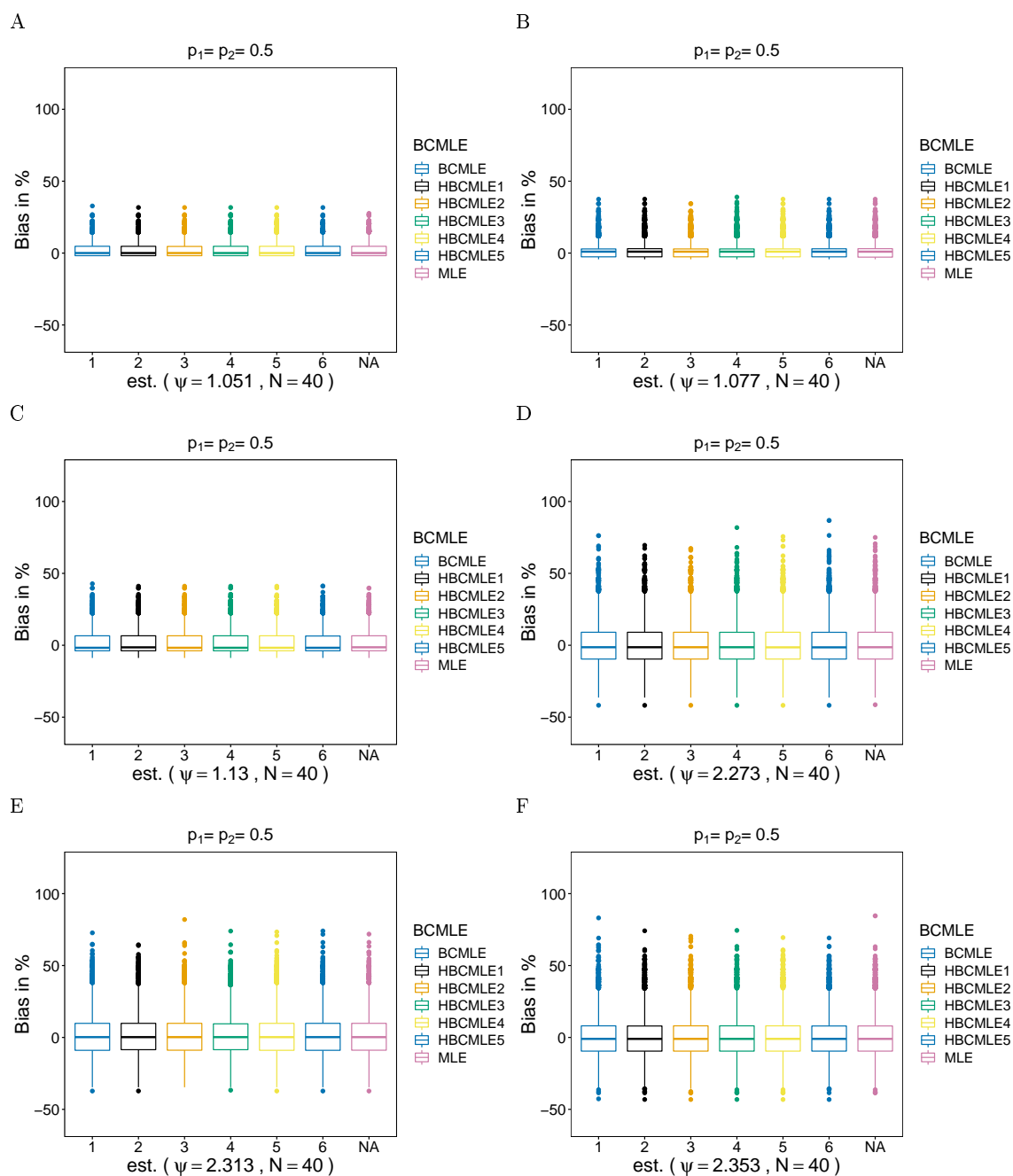


Figure 55: Same as Figure 51

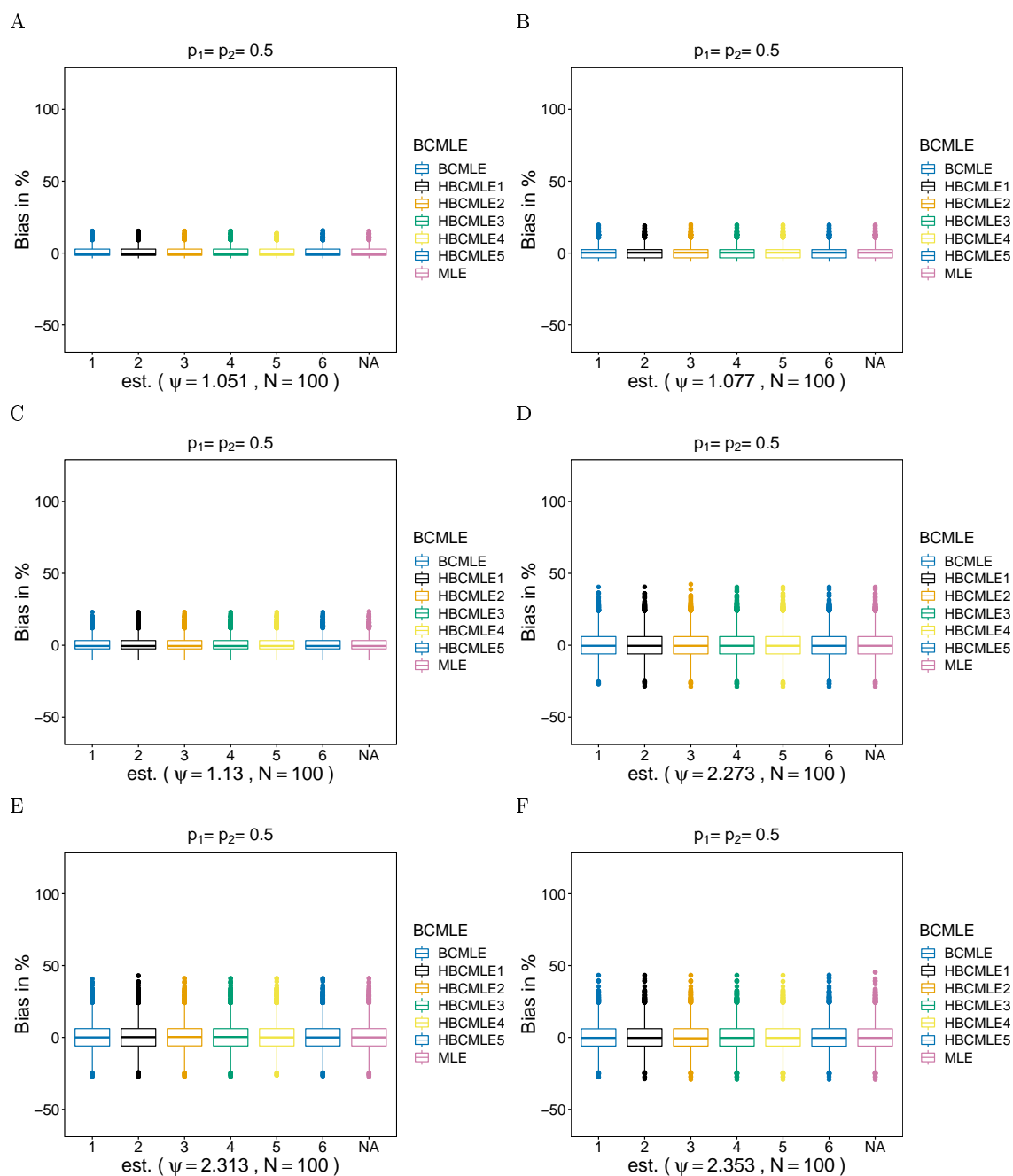


Figure 56: Same as Figure 51

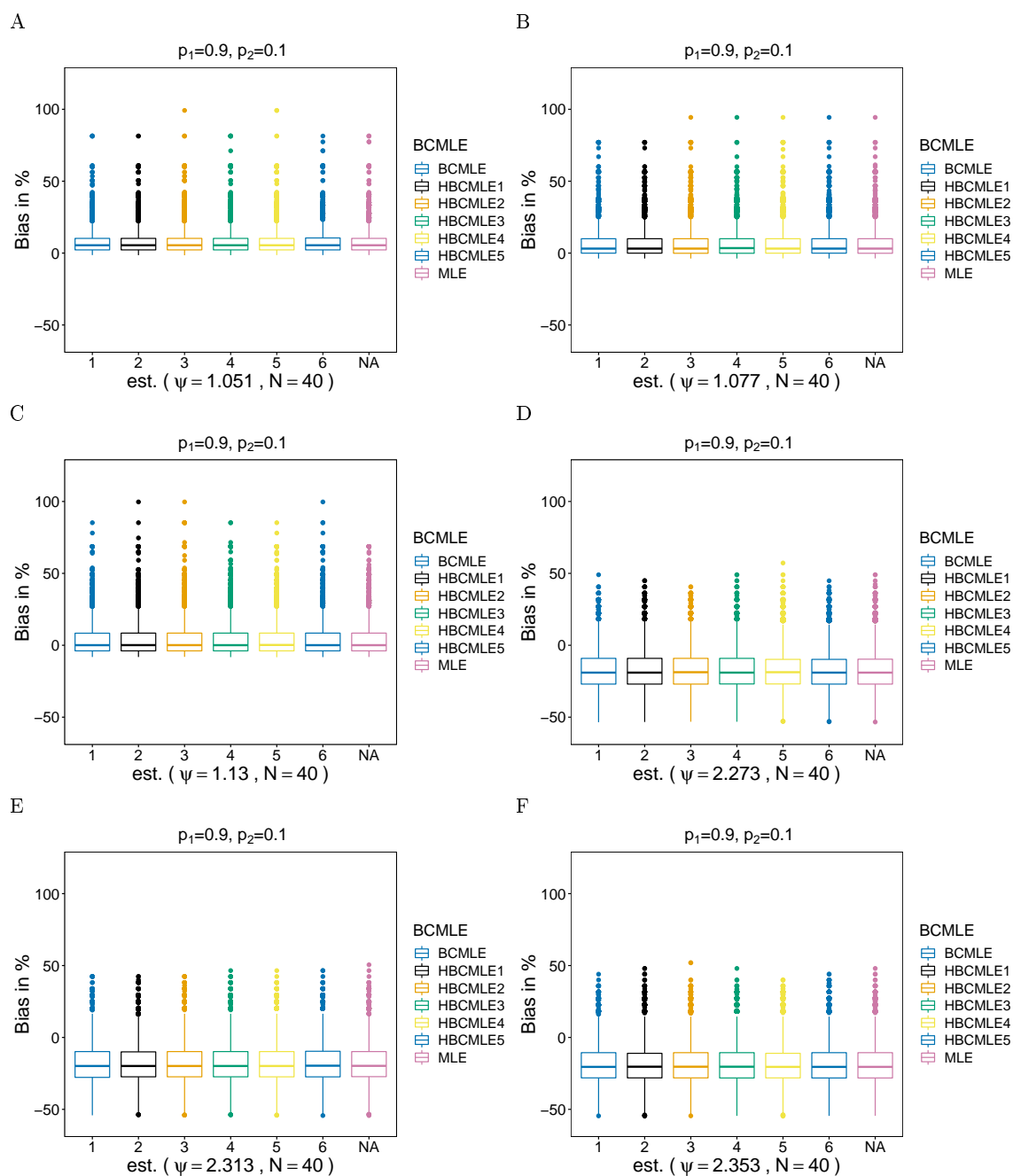


Figure 57: Same as Figure 51



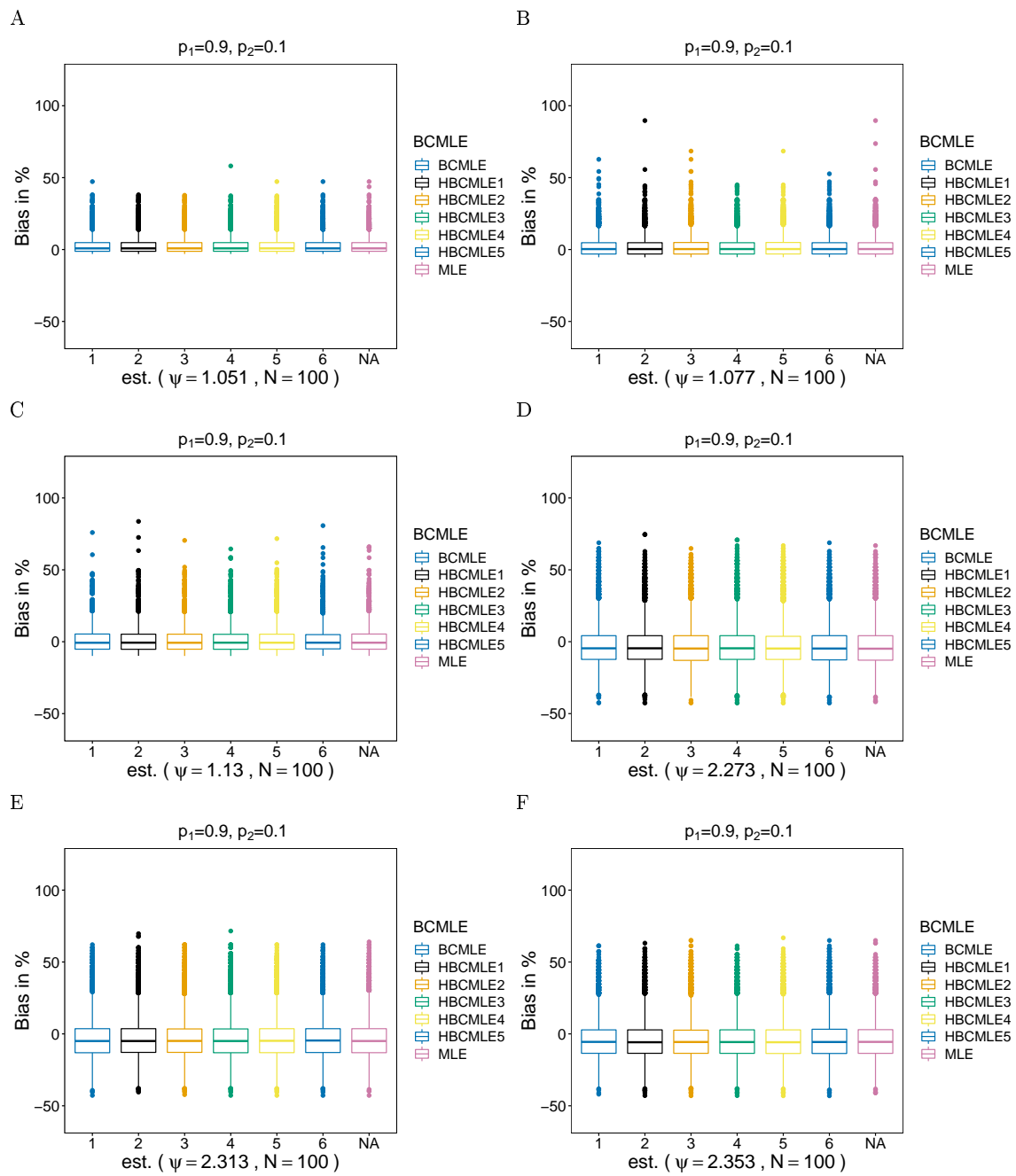


Figure 58: Same as Figure 51

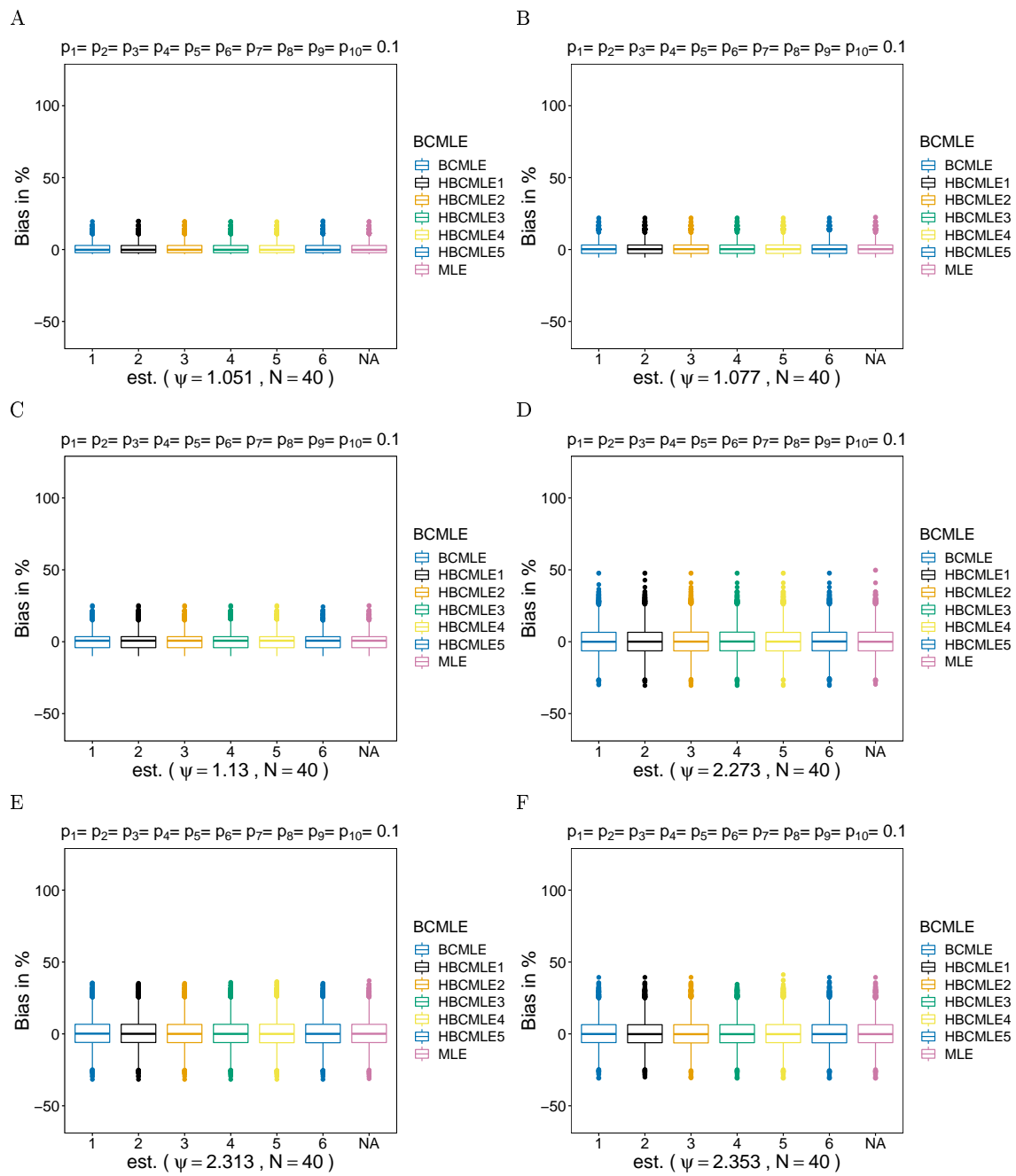


Figure 59: Same as Figure 51

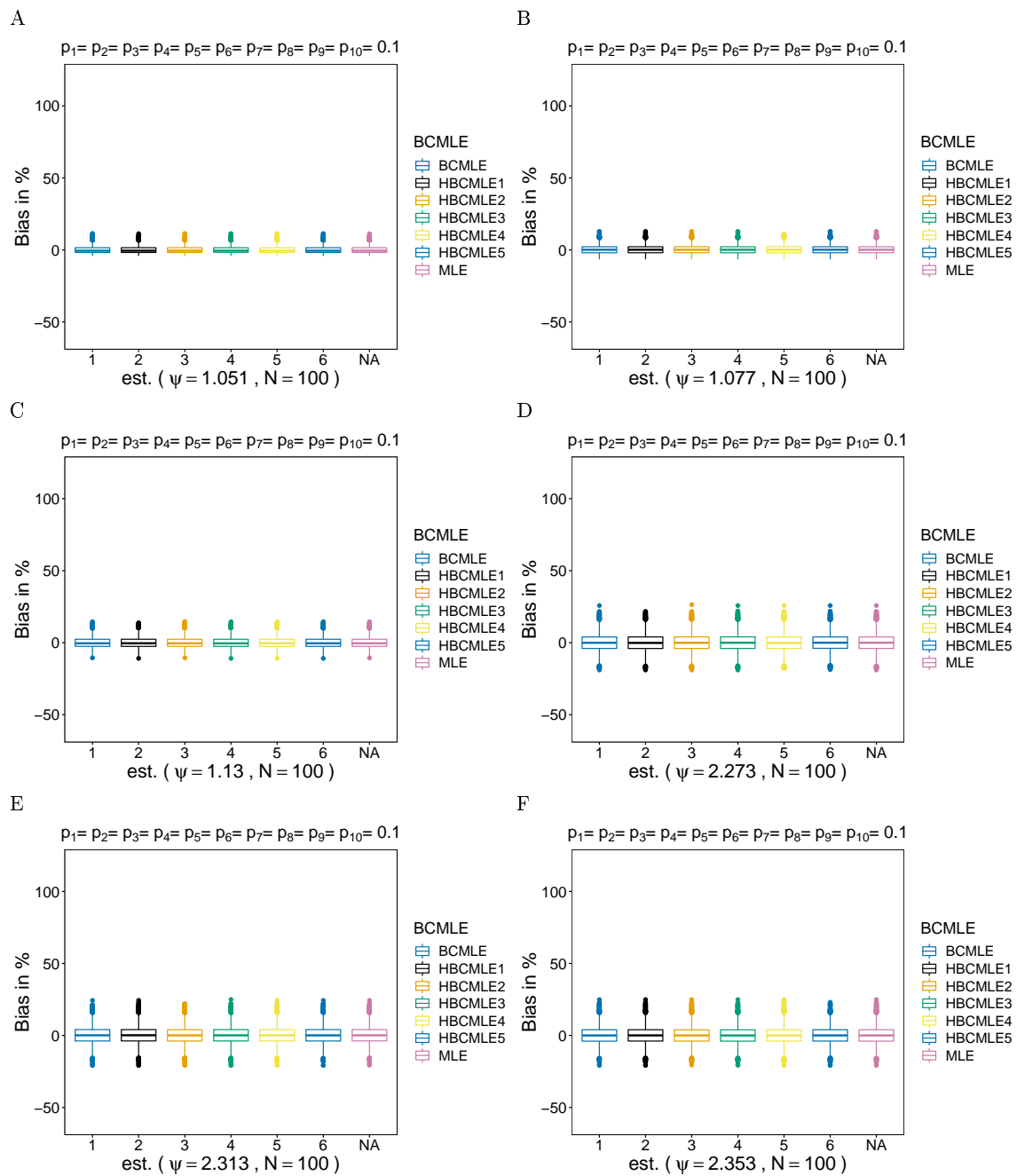


Figure 60: Same as Figure 51

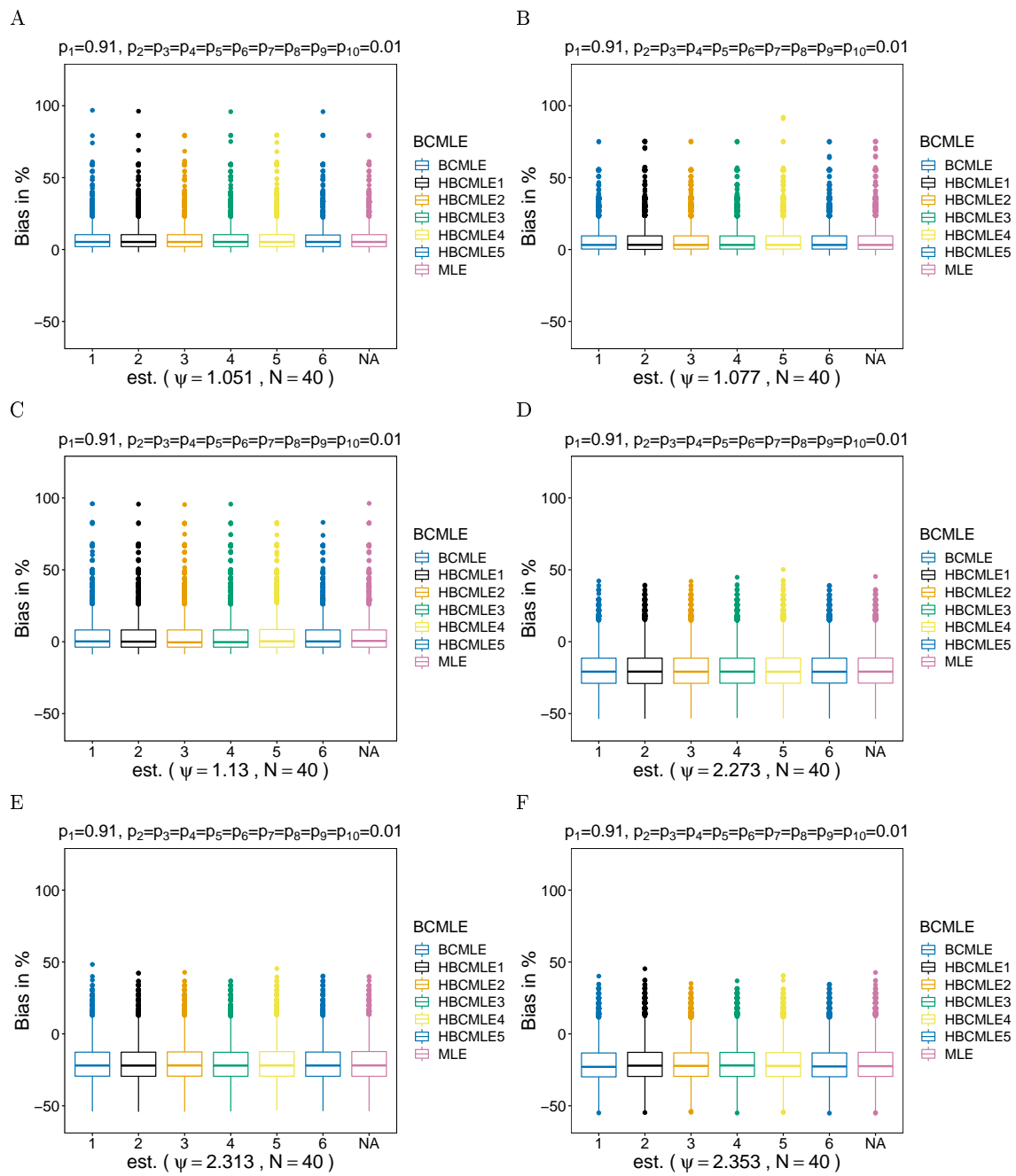


Figure 61: Same as Figure 51

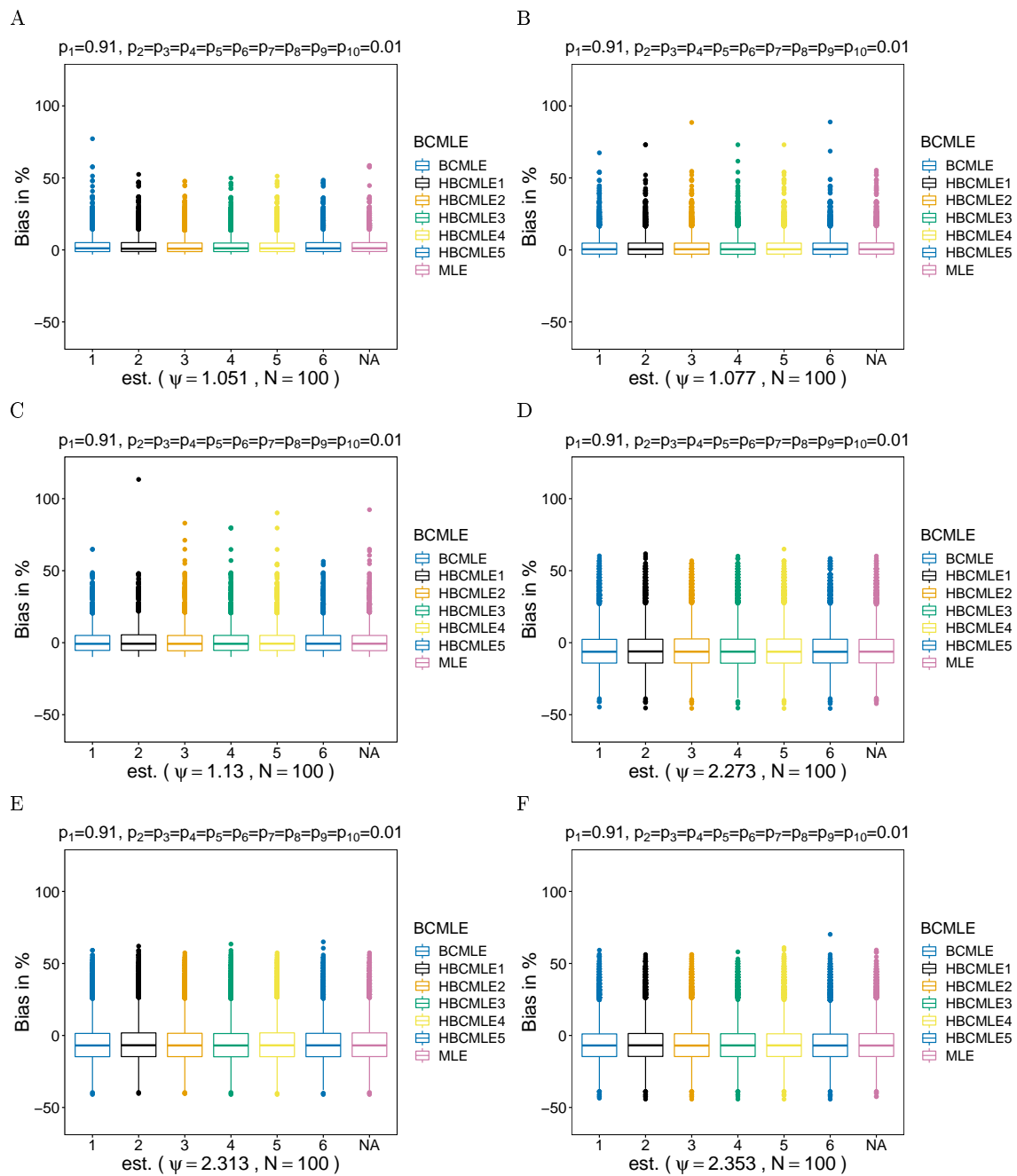


Figure 62: Same as Figure 51