

SUCCOTASH sims when α , the covariance, and the grid are known.

David Gerard

January 4, 2016

Abstract

This simulation study looks at the model assumed in the second step of SUCCOTASH. I compare SUCCOTASH with the second step of LEAPP. No other confounder adjustment procedure is applicable for comparison when assuming this model. I assume that α , Σ , and τ are all known.

1 Model Description

$$Y_{p \times 1} = \beta_{p \times 1} + \alpha_{p \times k} Z_{k \times 1} + E_{p \times 1}, \quad (1)$$

such that

- α is known.
- $E \sim N_p(0, I_p)$, so we explore homoscedastic case.

2 Procedure

- $p = 100$,
- $k \in \{5, 10, 50\}$,
- $\beta_j \sim N(0, \tau_k^2)$ w.p. π_k ,
- $\tau_k^2 = 0, 1, 100$ for $k = 0, 1, 2$ when we have a three mixture and $\tau_k = 0, 100$ for $k = 0, 1$ when we have a two mixture model,
- $\pi \in \{(0.5, 0.5), (0.9, 0.1), (1, 0, 0), (0.9, 0.1, 0), (0.9, 0, 0.1), (0.5, 0.5, 0), (0.5, 0, 0.5), (0.5, 0.25, 0.25)\}$
- $Z_j \stackrel{i.i.d.}{\sim} N(0, 1)$,
- $\alpha_{ij} \stackrel{i.i.d.}{\sim} N(0, 1)$,
- 400 iterations for each π by k combination, sampling a new Z and α at each iteration.
- I did not regularize the estimates of π for SUCCOTASH.
- At each iteration, I calculated the Sum of Squared Errors (SSE) for the posterior means under SUCCOTASH, and the estimates of β given by the second step of LEAPP.
- I also calculated the SSE when using just Y to estimate β (called OLS in Figures and Tables below).
- I also calculated $\hat{\pi}_0$ given by SUCCOTASH and LEAPP at each iteration.
- LEAPP uses an L_1 penalty, so I called its $\hat{\pi}_0$ to just be the proportion of elements of β it sets to 0.

- The only comparable procedure using Model (1) is the second step of LEAPP [Sun et al., 2012].
- CATE [Wang et al., 2015] is not applicable because its model for its second step is

$$Y_{p \times 1} \sim N_p(\beta_{p \times 1} + \alpha_{p \times k} \gamma_{k \times 1}, \alpha \alpha^T + I_p), \quad (2)$$

where γ describes the linear relationship between the observed and unobserved variables.

- RUV [Gagnon-Bartsch et al., 2013] is not applicable because we don't assume we have any control genes.
- SVA [Leek and Storey, 2008] is not applicable because it doesn't use this two-step procedure.

3 Results

SUCCOTASH always beats LEAPP in terms of SSE (Table 1), especially when k is large. LEAPP seems to perform horribly whenever there are a lot of confounders.

When k is small ($k = 5$ or 10) and we don't include $\tau_1 = 1$, SUCCOTASH estimates π_0 fairly accurately (Table 3). This accuracy decreases when k is large ($k = 50$). But increasing p corrects for this (Table 4).

4 SSE Boxplots

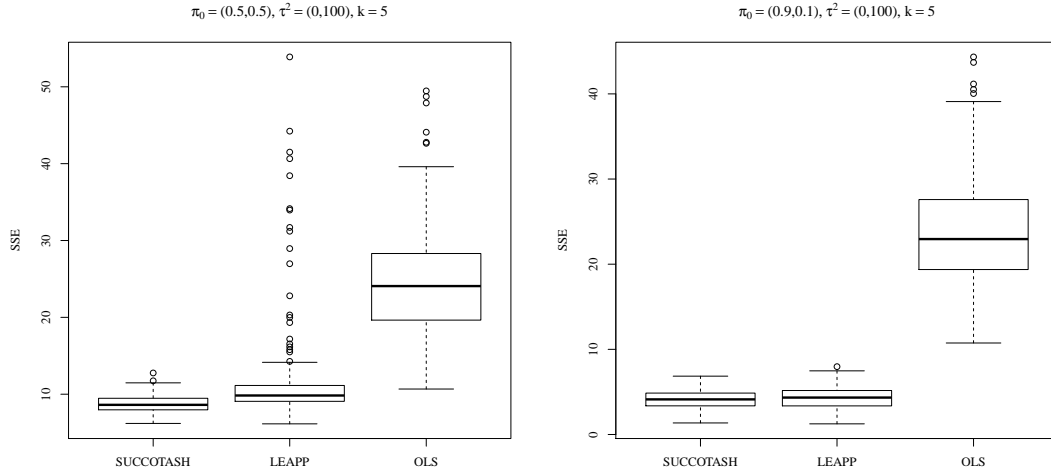


Table 1: Average Sum of Squared Errors for SUCCOTASH, LEAPP, and OLS at given k and π values.

| k | π_0 | π_1 | π_2 | SUCCOTASH | LEAPP | OLS |
|-----|---------|---------|---------|-----------|-------|------|
| 5 | .5 | NA | .5 | 8.7 | 10.9 | 24.4 |
| 5 | .9 | NA | .1 | 4.1 | 4.3 | 23.6 |
| 5 | 1 | 0 | 0 | 0.2 | 0.8 | 23.7 |
| 5 | .9 | .1 | 0 | 3.0 | 3.5 | 23.7 |
| 5 | .9 | 0 | .1 | 4.2 | 4.5 | 23.6 |
| 5 | .5 | .5 | 0 | 5.9 | 7.3 | 23.7 |
| 5 | .5 | 0 | .5 | 8.9 | 10.7 | 24.2 |
| 5 | .5 | .25 | .25 | 7.6 | 8.4 | 23.7 |
| 10 | .5 | NA | .5 | 10.0 | 19.4 | 32.1 |
| 10 | .9 | NA | .1 | 4.3 | 4.6 | 32.5 |
| 10 | 1 | 0 | 0 | 0.2 | 1.3 | 32.4 |
| 10 | .9 | .1 | 0 | 3.0 | 3.7 | 32.7 |
| 10 | .9 | 0 | .1 | 4.3 | 4.7 | 31.9 |
| 10 | .5 | .5 | 0 | 6.0 | 7.4 | 32.7 |
| 10 | .5 | 0 | .5 | 10.1 | 18.2 | 32.5 |
| 10 | .5 | .25 | .25 | 8.1 | 9.7 | 32.5 |
| 50 | .5 | NA | .5 | 54.2 | 99.5 | 70.8 |
| 50 | .9 | NA | .1 | 7.5 | 98.8 | 70.9 |
| 50 | 1 | 0 | 0 | 0.0 | 98.5 | 69.8 |
| 50 | .9 | .1 | 0 | 3.1 | 98.7 | 71.1 |
| 50 | .9 | 0 | .1 | 7.4 | 97.9 | 71.2 |
| 50 | .5 | .5 | 0 | 7.4 | 99.3 | 70.9 |
| 50 | .5 | 0 | .5 | 52.8 | 100.4 | 71.2 |
| 50 | .5 | .25 | .25 | 23.2 | 99.9 | 71.1 |

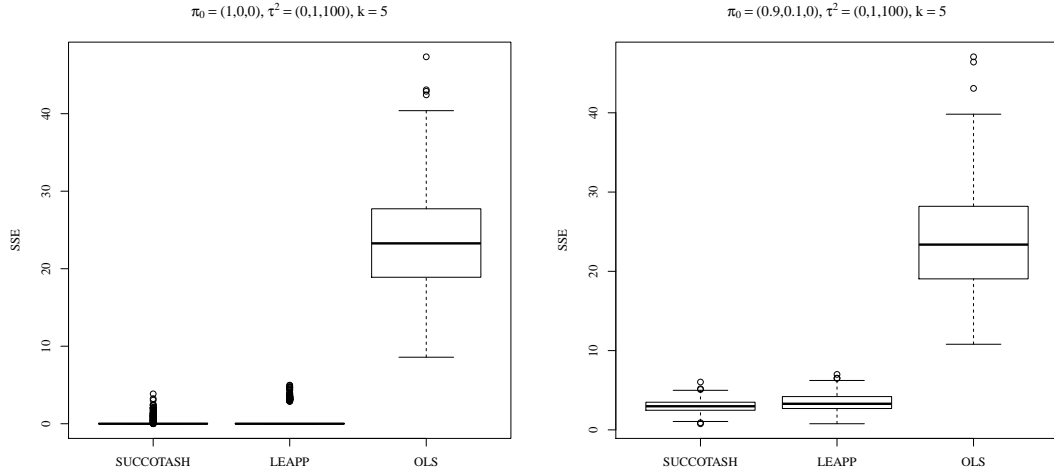


Table 2: Average Sum of Squared Errors for SUCCOTASH, LEAPP, and OLS at given k and π values. Here, $p = 500$.

| k | π_0 | π_1 | π_2 | SUCCOTASH | LEAPP | OLS |
|-----|---------|---------|---------|-----------|-------|-------|
| 50 | .5 | NA | .5 | 22.1 | 134.4 | 159.0 |
| 50 | .9 | NA | .1 | 10.0 | 10.9 | 159.1 |

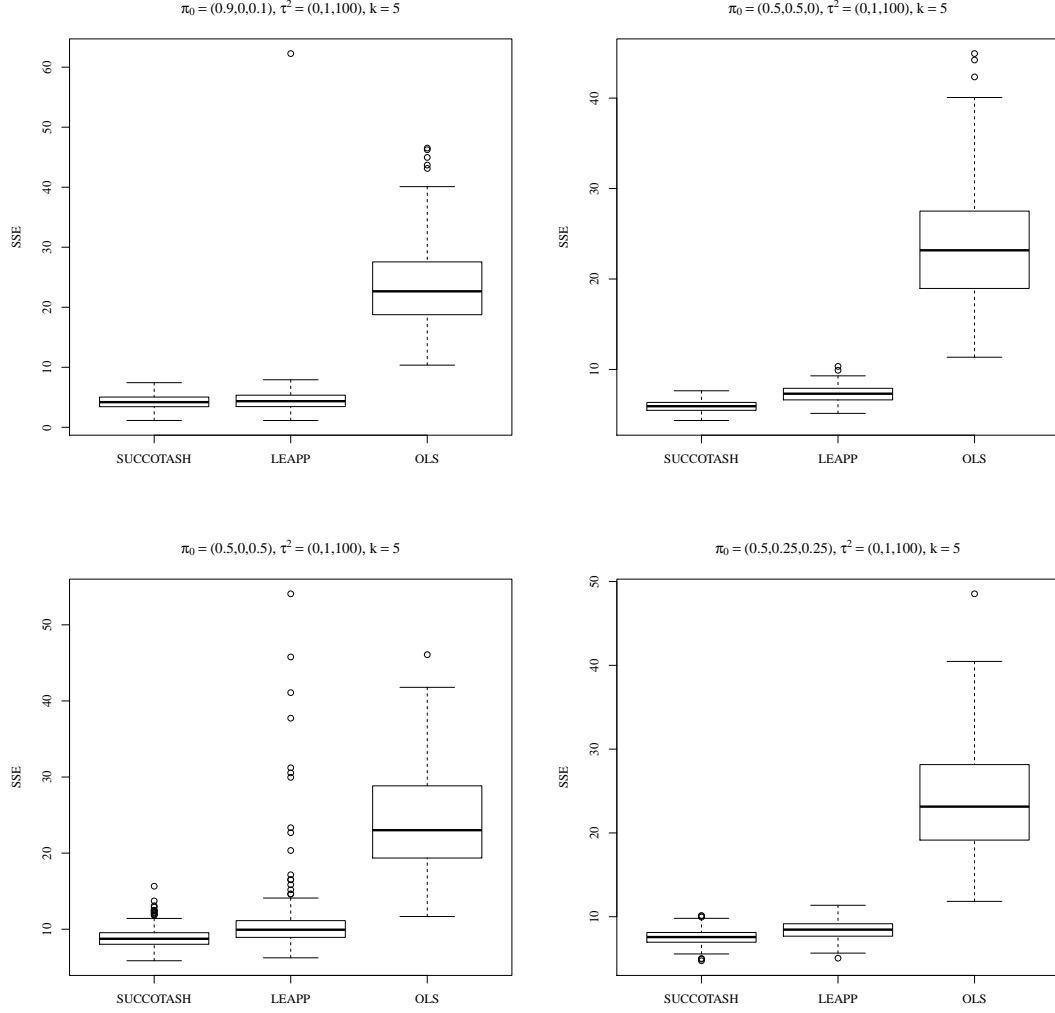


Table 3: Mean $\hat{\pi}_0$ for SUCCOTASH and LEAPP at given k and π values.

| k | π_0 | π_1 | π_2 | SUCCOTASH | LEAPP |
|-----|---------|---------|---------|-----------|-------|
| 5 | .5 | NA | .5 | .51 | .62 |
| 5 | .9 | NA | .1 | .90 | .92 |
| 5 | 1 | 0 | 0 | .98 | 1.00 |
| 5 | .9 | .1 | 0 | .92 | .99 |
| 5 | .9 | 0 | .1 | .86 | .92 |
| 5 | .5 | .5 | 0 | .57 | .98 |
| 5 | .5 | 0 | .5 | .48 | .62 |
| 5 | .5 | .25 | .25 | .57 | .80 |
| 10 | .5 | NA | .5 | .52 | .58 |
| 10 | .9 | NA | .1 | .91 | .92 |
| 10 | 1 | 0 | 0 | .98 | 1.00 |
| 10 | .9 | .1 | 0 | .94 | .99 |
| 10 | .9 | 0 | .1 | .88 | .92 |
| 10 | .5 | .5 | 0 | .65 | .98 |
| 10 | .5 | 0 | .5 | .50 | .58 |
| 10 | .5 | .25 | .25 | .61 | .79 |
| 50 | .5 | NA | .5 | .66 | .17 |
| 50 | .9 | NA | .1 | .92 | .20 |
| 50 | 1 | 0 | 0 | 1.00 | .20 |
| 50 | .9 | .1 | 0 | 1.00 | .21 |
| 50 | .9 | 0 | .1 | .93 | .20 |
| 50 | .5 | .5 | 0 | 1.00 | .20 |
| 50 | .5 | 0 | .5 | .68 | .17 |
| 50 | .5 | .25 | .25 | .81 | .19 |

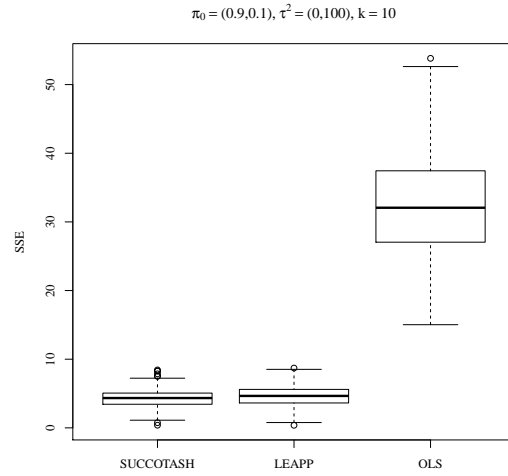
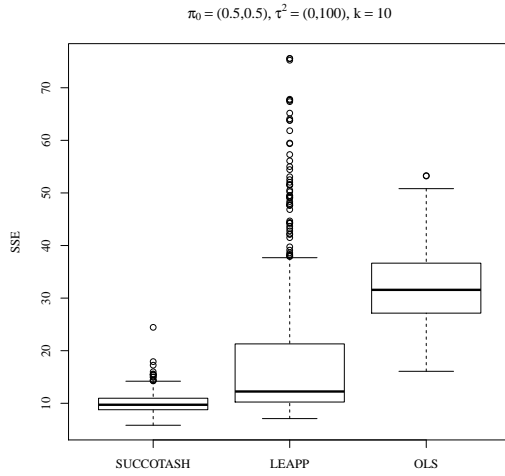
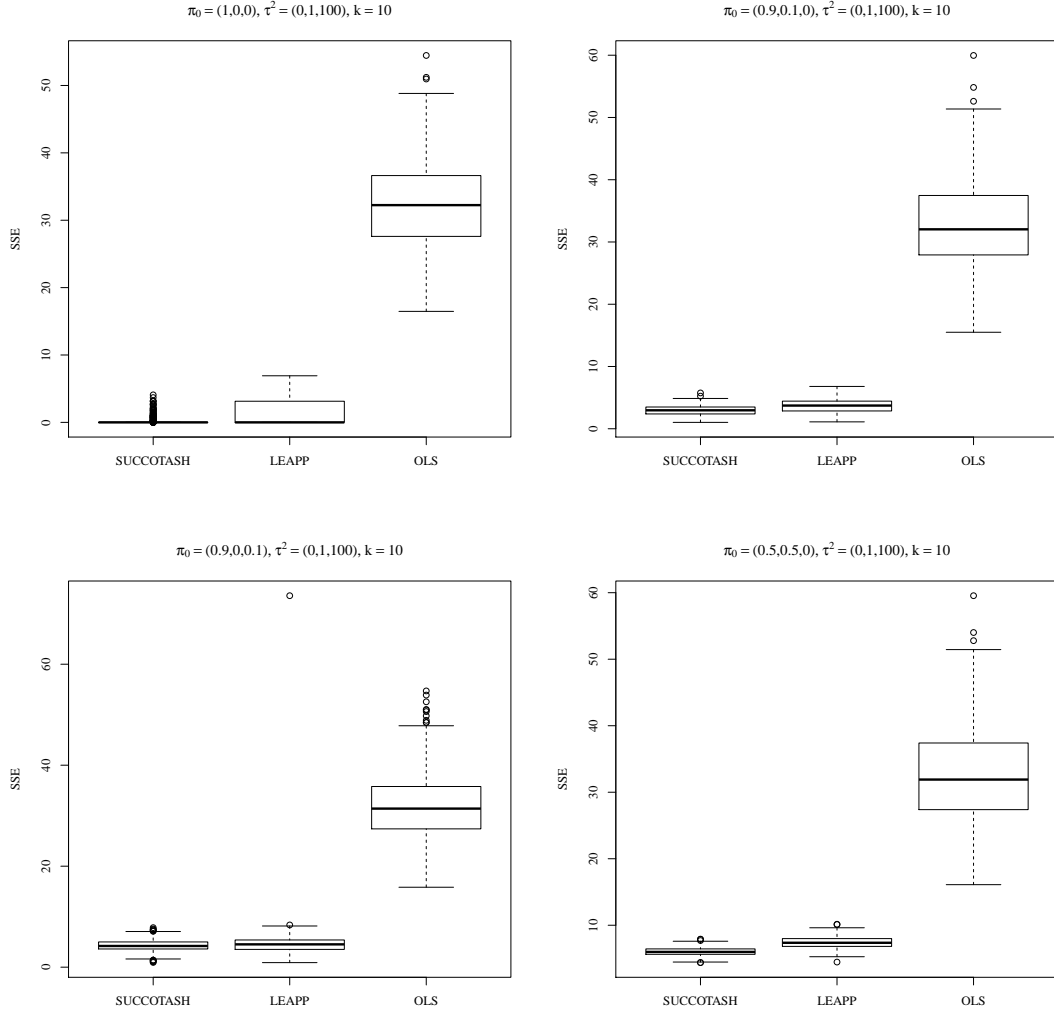
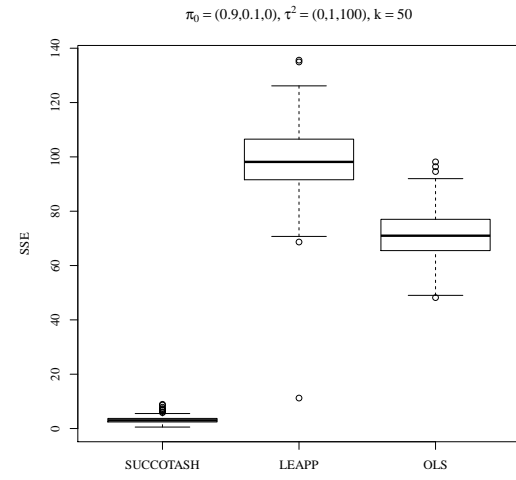
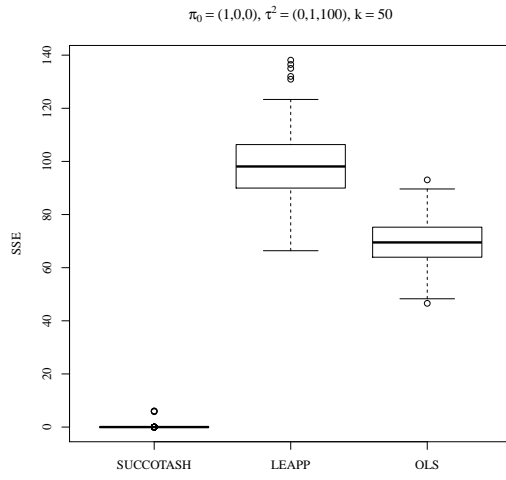
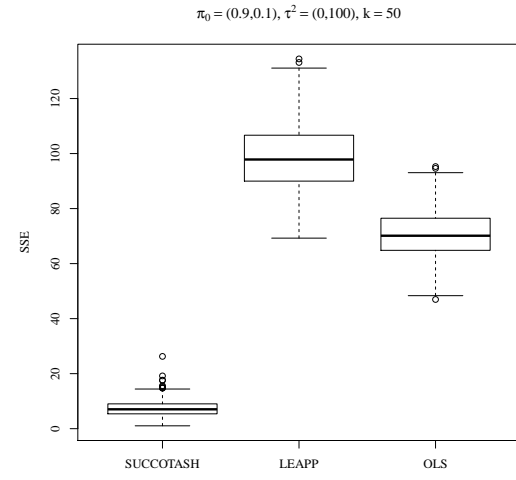
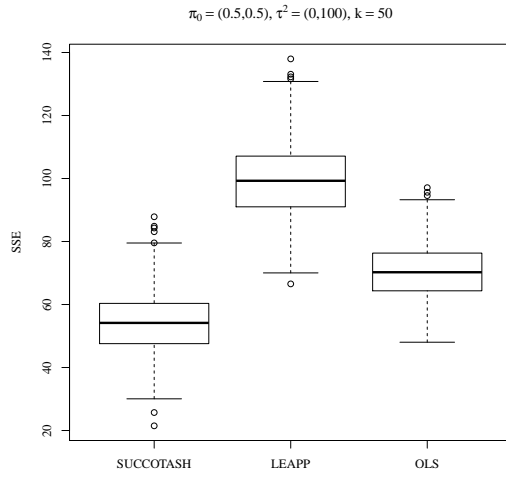
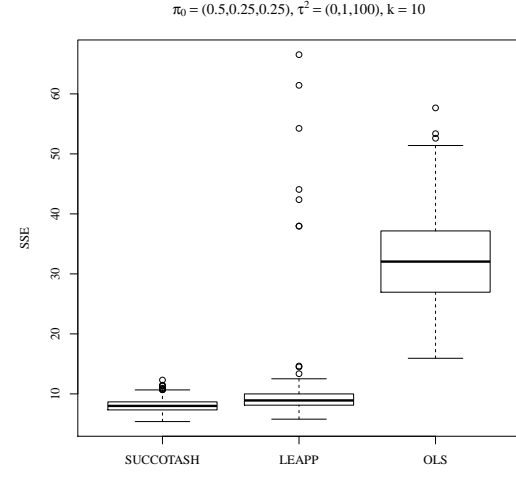
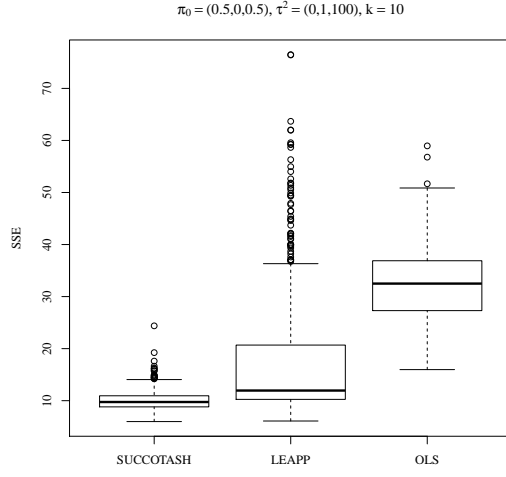
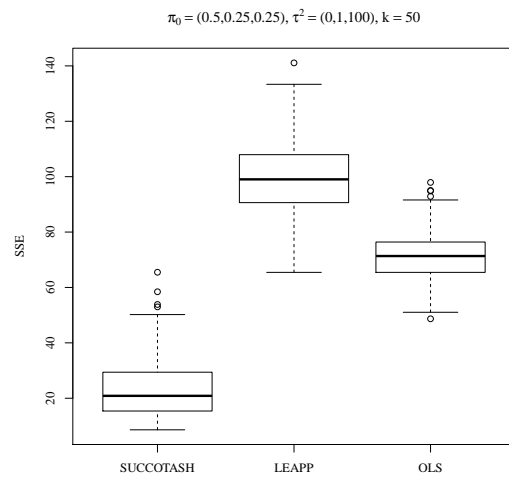
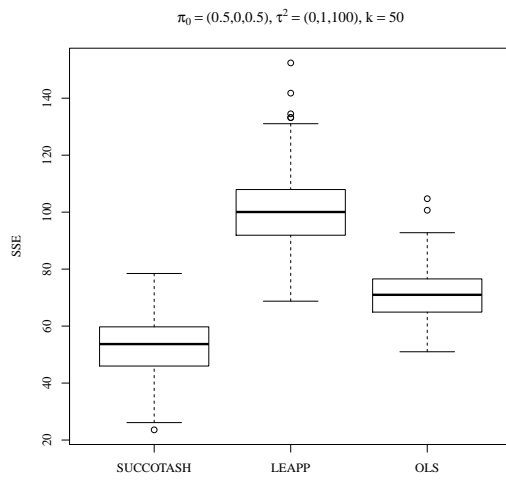
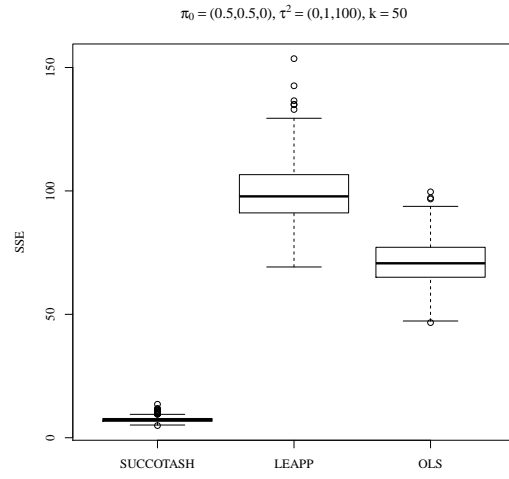
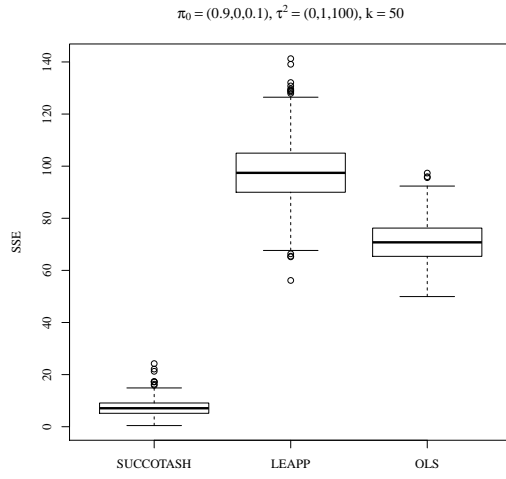


Table 4: Mean $\hat{\pi}_0$ for SUCCOTASH and LEAPP at given k and π values. Here, $p = 500$.

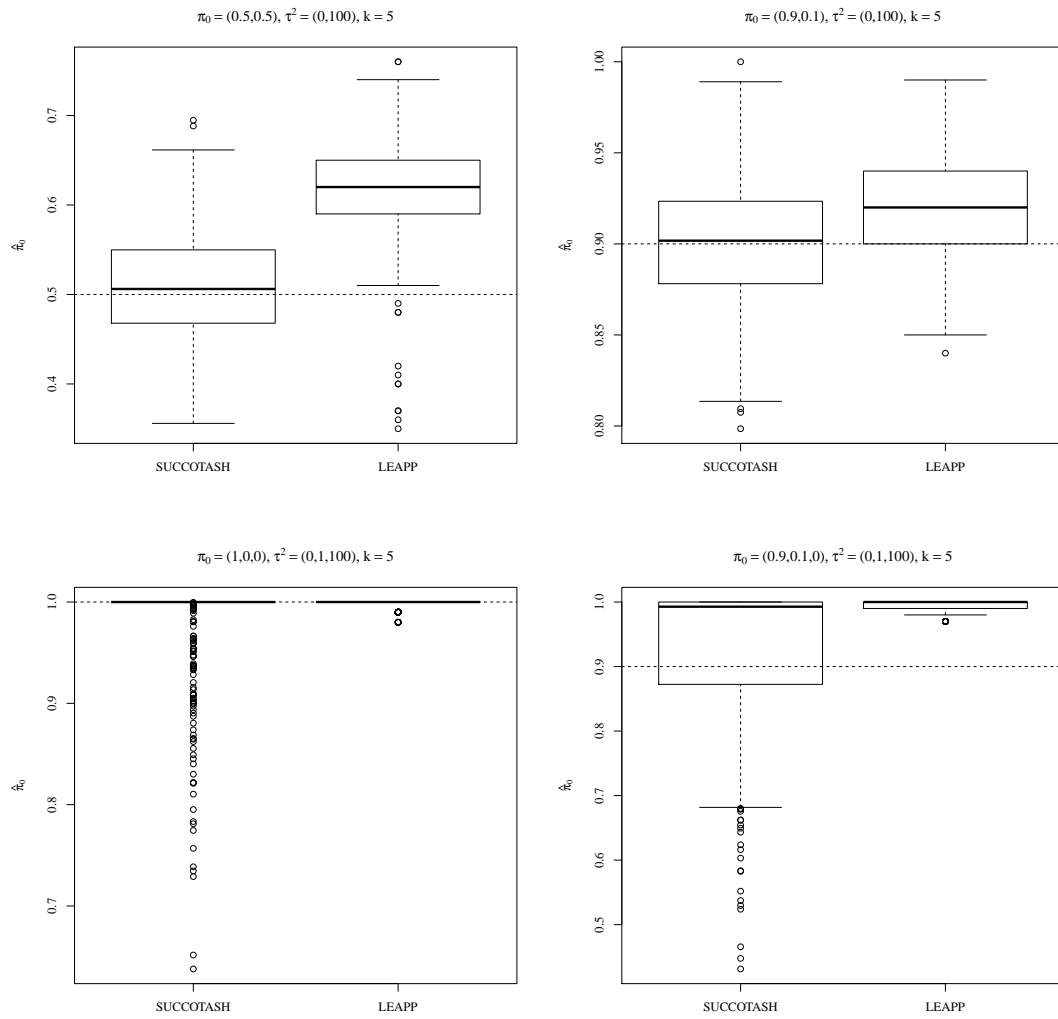
| k | π_0 | π_1 | π_2 | SUCCOTASH | LEAPP |
|-----|---------|---------|---------|-----------|-------|
| 50 | .5 | NA | .5 | .52 | .45 |
| 50 | .9 | NA | .1 | .90 | .93 |

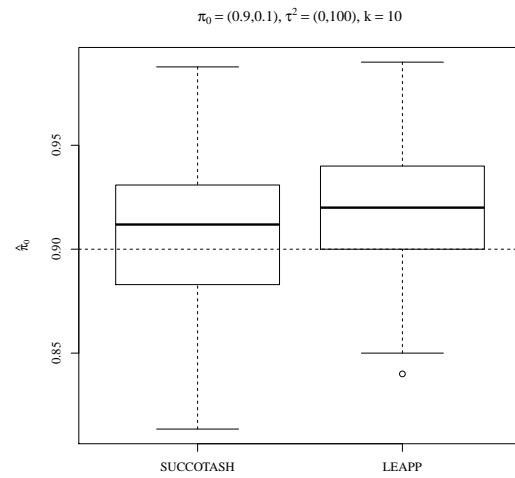
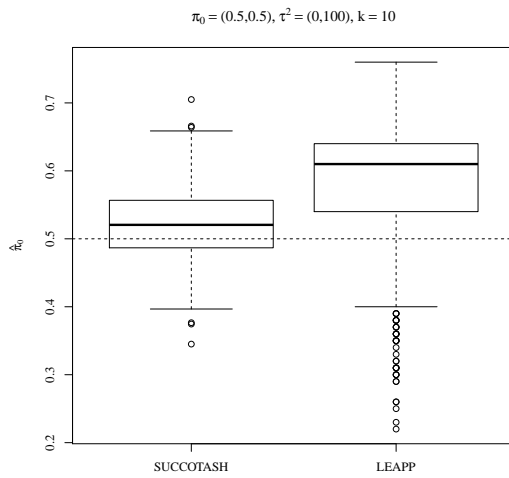
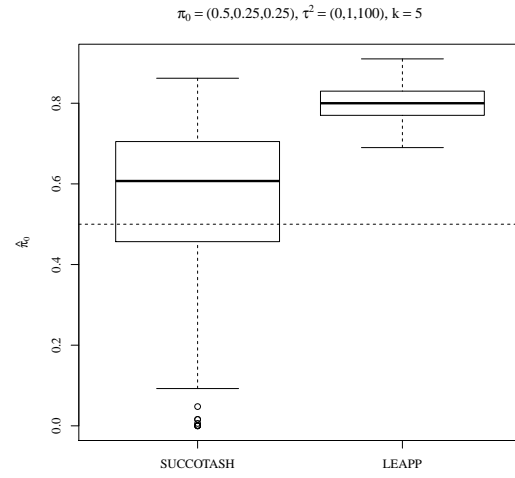
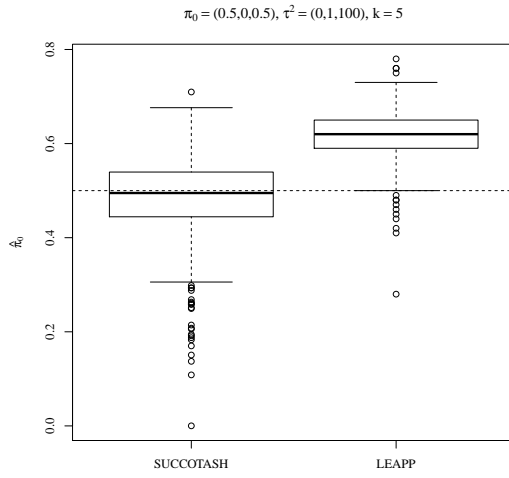
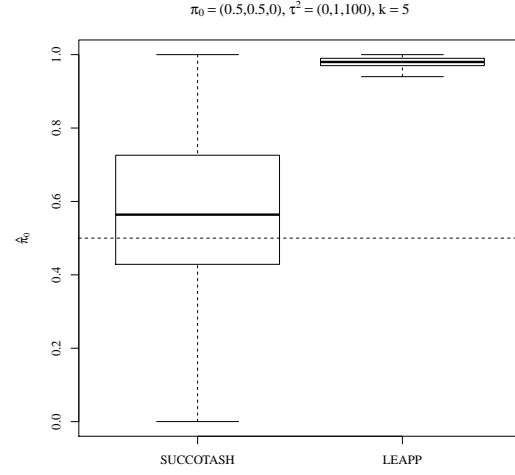
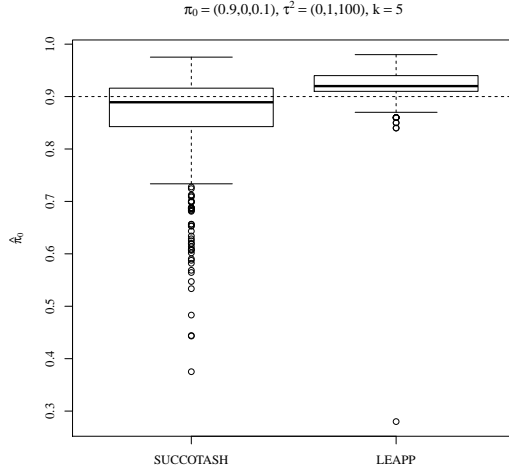


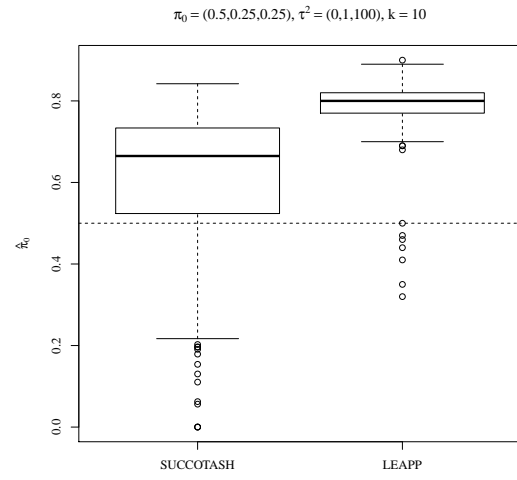
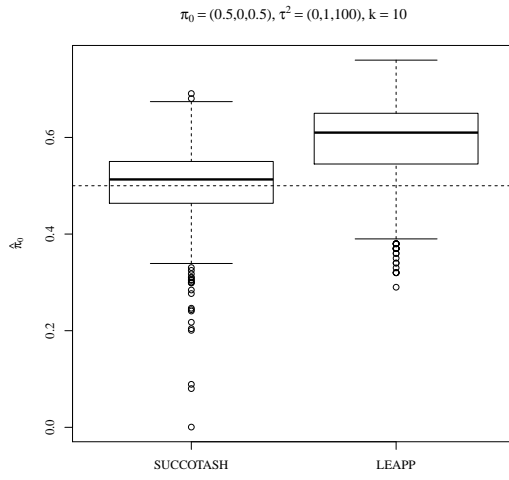
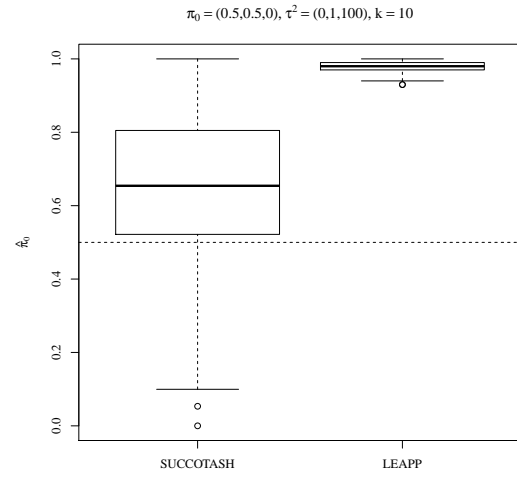
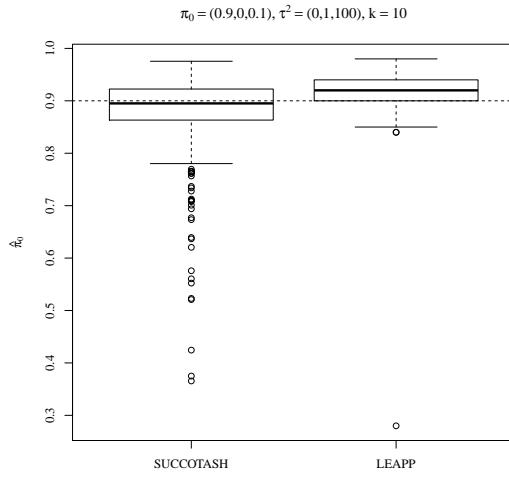
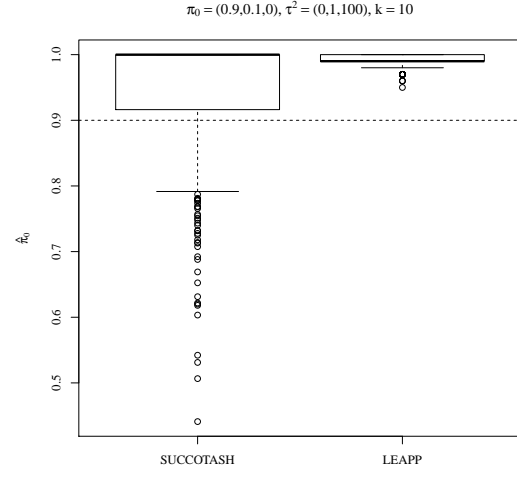
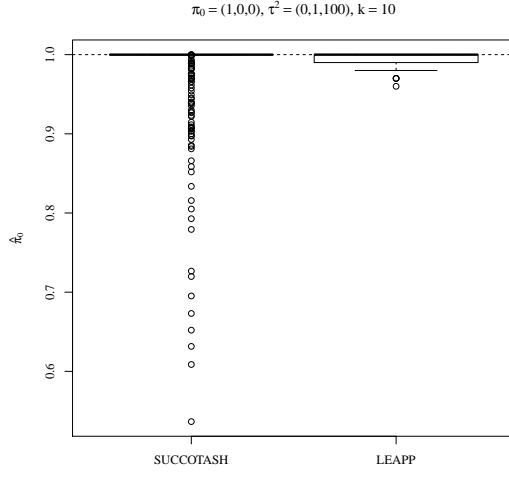


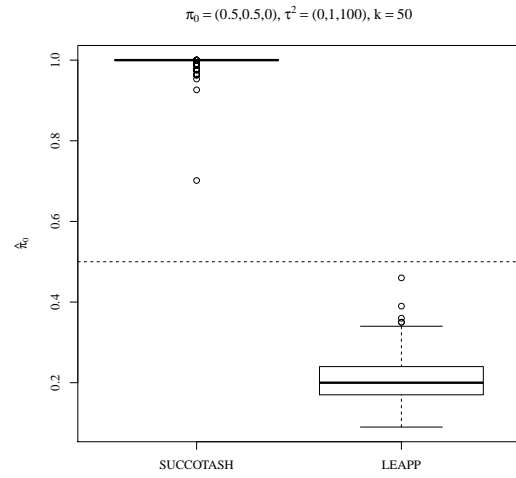
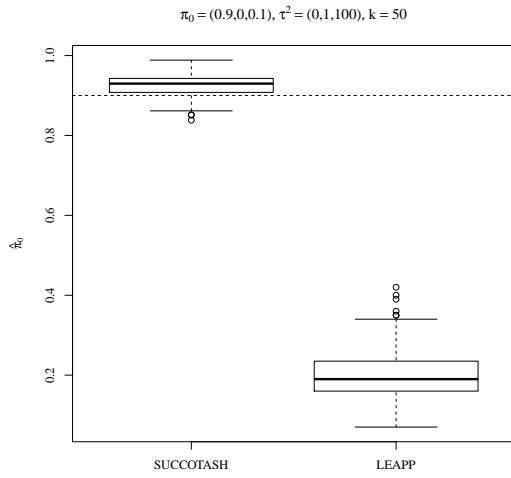
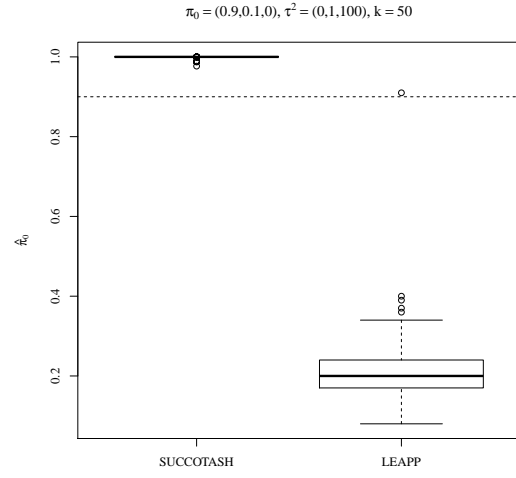
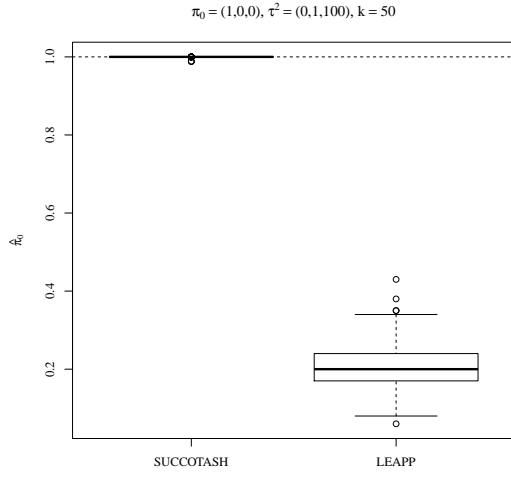
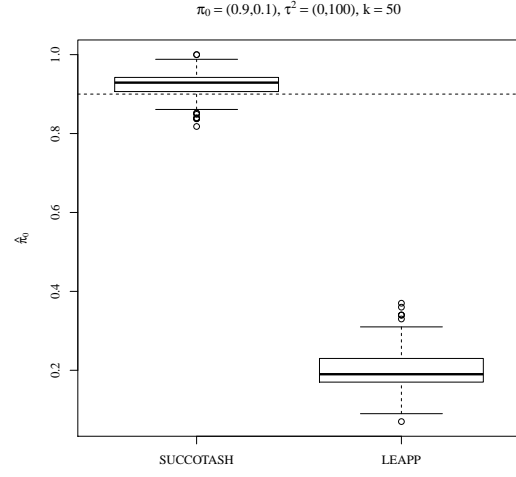
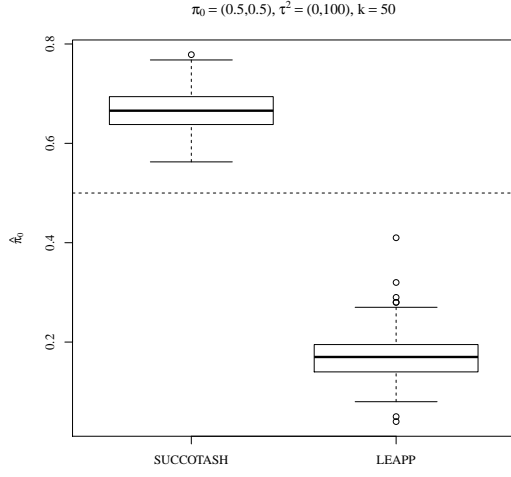


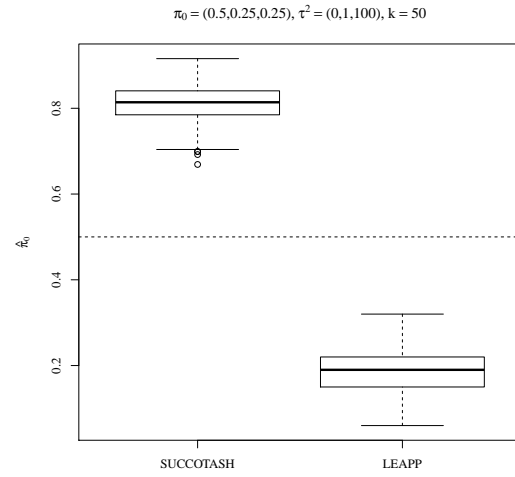
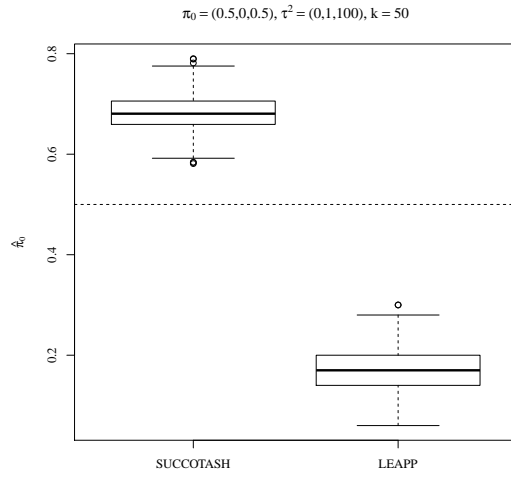
5 $\hat{\pi}_0$ Boxplots



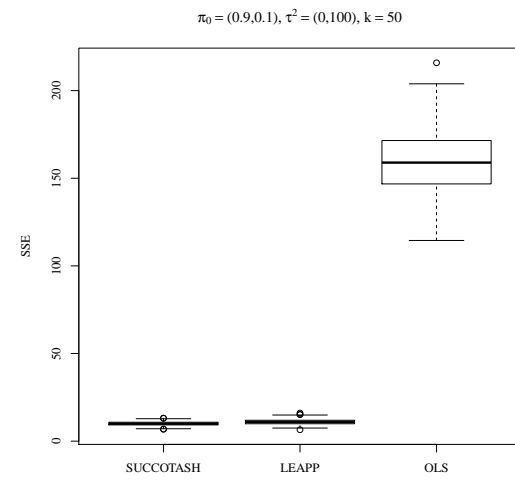
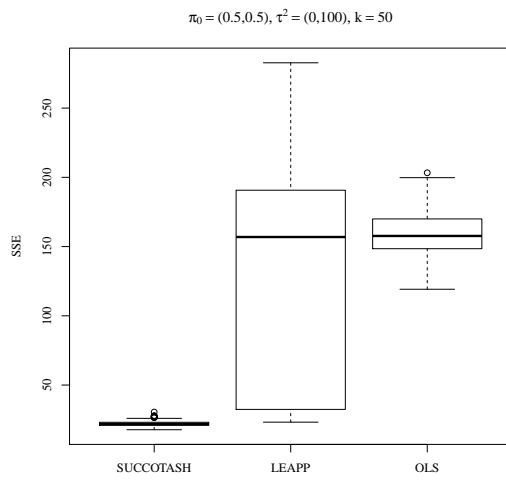


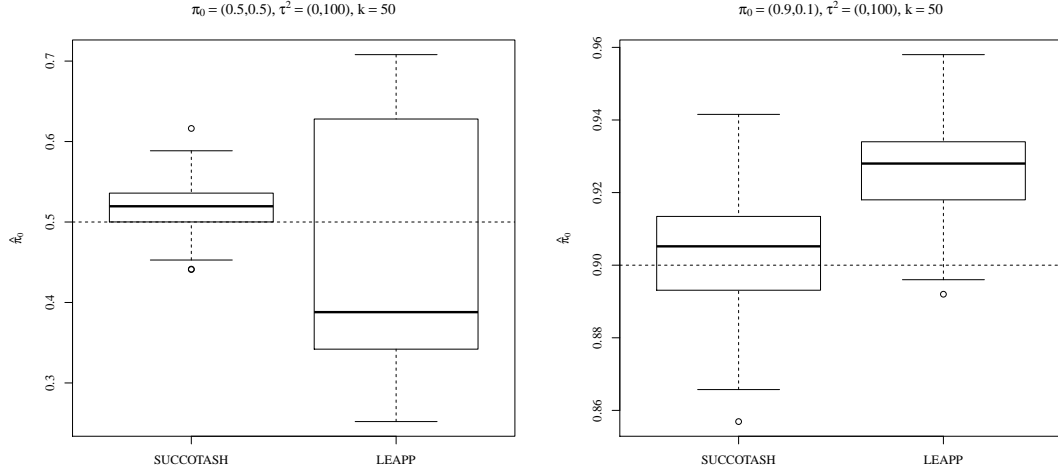






6 Boxplots when $p = 500$





References

- J Gagnon-Bartsch, L Jacob, and TP Speed. Removing unwanted variation from high dimensional data with negative controls. Technical report, Technical Report 820, Department of Statistics, University of California, Berkeley, 2013.
- Jeffrey T Leek and John D Storey. A general framework for multiple testing dependence. *Proceedings of the National Academy of Sciences*, 105(48):18718–18723, 2008.
- Yunting Sun, Nancy R Zhang, Art B Owen, et al. Multiple hypothesis testing adjusted for latent variables, with an application to the agemap gene expression data. *The Annals of Applied Statistics*, 6(4):1664–1688, 2012.
- Jingshu Wang, Qingyuan Zhao, Trevor Hastie, and Art B Owen. Confounder adjustment in multiple hypotheses testing. *arXiv preprint arXiv:1508.04178*, 2015.