

# Supplemental material on “Improving the computational performance of TCLUST through ensemble initialization”

Pedro C. Álvarez-Esteban, Luis A. García-Escudero,  
Agustín Mayo-Iscar, Javier Crespo-Guerrero

## 1 R Code for simula\_tclust

```
1 simula_tclust <- function(n, p = 4, G = 3, scenario = 2, balanced = 1, symmetry = 1)
2 {
3   if (G != 3 & G != 6)
4     stop("Invalid G")
5   if (p < 2)
6     stop("Invalid p")
7   if (symmetry != 1 & symmetry != 2)
8     stop("Invalid symmetry")
9   if (G == 6) {
10     nn <- n/2
11   }
12   else {
13     nn <- n
14   }
15   if (balanced == 1) {
16     nn1 = nn2 = nn3 = ceiling(nn * 0.3)
17   }
18   else {
19     nn1 = ceiling(nn * 0.25)
20     nn2 = ceiling(nn * 0.3)
21     nn3 = ceiling(nn * 0.35)
22   }
23   if (symmetry == 1 ){
24     cc1 <- -20
25     cc2 <- -50
26   } else {
27     cc1 <- 0
28     cc2 <- 0
29   }
30 }
31
32 nh <- nn1 + nn2 + nn3
33 Y <- matrix(NA, ncol = p, nrow = n)
34 if (scenario == 1) {
35   if (G == 3) {
36     Y[1:nn1, ] <- MASS::mvrnorm(nn1, c(4, 4), matrix(c(1, 0, 0, 1), nrow = 2))
37     Y[(nn1 + 1):(nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(-4, 4), matrix(c(1, 0, 0, 1), nrow = 2))
38     Y[(nn1 + nn2 + 1):(nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(0, 0), matrix(c(1, 0, 0, 1),
39     nrow = 2))
40     Y[(nh + 1):n, ] <- matrix(runif(p * (n - nh), cc1, 20), ncol = p)
41     if (p > 2) {
42       Y[1:n, 3:p] <- MASS::mvrnorm(n, rep(0, p - 2), diag(p - 2))
43     }
44   }
45   else {
46     Y[1:nn1, ] <- MASS::mvrnorm(nn1, c(4, 4), matrix(c(1, 0, 0, 1), nrow = 2))
47     Y[(nn1 + 1):(nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(-4, 4), matrix(c(1, 0, 0, 1), nrow = 2))
     Y[(nn1 + nn2 + 1):(nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(0, 0), matrix(c(1, 0, 0, 1),
      nrow = 2))
48     Y[(nh + 1):(nh + nn1), ] <- MASS::mvrnorm(nn1, c(24, 0), matrix(c(1, 0, 0, 1), nrow = 2))
49     Y[(nh + nn1 + 1):(nh + nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(16, 0), matrix(c(1, 0, 0, 1),
      nrow = 2))
50     Y[(nh + nn1 + nn2 + 1):(nh + nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(20, 4),
      matrix(c(1, 0, 0, 1), nrow = 2))
51     Y[(2 * nh + 1):n, ] <- matrix(runif((n - 2 * nh) * p, cc1, 40), ncol = p)
```

```

52   if (p > 2) {
53     Y[1:n, 3:p] <- MASS::mvrnorm(n, rep(0, p - 2), diag(p - 2))
54   }
55 }
56 if (scenario == 2) {
57   alp1 <- runif(1, 0, 2 * pi)
58   u1 <- matrix(c(cos(alp1), sin(alp1), -sin(alp1), cos(alp1)), ncol = 2)
59   alp2 <- runif(1, 0, 2 * pi)
60   u2 <- matrix(c(cos(alp2), sin(alp2), -sin(alp2), cos(alp2)), ncol = 2)
61   if (G == 3) {
62     Y[1:nn1, ] <- MASS::mvrnorm(nn1, c(20, 20), t(u1) %*% diag(c(1, 9^2)) %*% u1)
63     Y[(nn1 + 1):(nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(-20, -20), t(u1) %*% diag(c(9^2, 1)) %*%
64       u1)
65     Y[(nn1 + nn2 + 1):(nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(0, 0), t(u1) %*%
66       diag(c(3^2, 3^2)) %*% u1)
67     Y[(nn1 + nn2 + nn3 + 1):n, ] <- matrix(runif(p * (n - nh), cc2, 50), ncol = p)
68     if (p > 2) {
69       Y[1:n, 3:p] <- MASS::mvrnorm(n, rep(0, p - 2), diag(p - 2))
70     }
71   } else {
72     Y[1:nn1, ] <- MASS::mvrnorm(nn1, c(20, 20), t(u1) %*% diag(c(1, 9^2)) %*% u1)
73     Y[(nn1 + 1):(nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(-20, -20), t(u1) %*% diag(c(9^2, 1)) %*%
74       u1)
75     Y[(nn1 + nn2 + 1):(nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(0, 0), t(u1) %*% diag(c(3^2,
76       3^2)) %*% u1)
77     Y[(nh + 1):(nh + nn1), ] <- MASS::mvrnorm(nn1, c(40, 20), t(u2) %*% diag(c(1, 9^2)) %*% u2)
78     Y[(nh + nn1 + 1):(nh + nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(0, -20), t(u2) %*% diag(c(9^2,
79       1)) %*% u2)
80     Y[(nh + nn1 + nn2 + 1):(nh + nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(20, 0), t(u2) %*%
81       diag(c(3^2, 3^2)) %*% u2)
82     Y[(2 * nh + 1):n, ] <- matrix(runif((n - 2 * nh) * p, cc2-10, 80), ncol = p)
83     if (p > 2) {
84       Y[1:n, 3:p] <- MASS::mvrnorm(n, rep(0, p - 2), diag(p - 2))
85     }
86   }
87 }
88 if (scenario == 3) {
89   if (G == 3) {
90     Y[1:nn1, ] <- MASS::mvrnorm(nn1, c(5, 5), matrix(c(1, 0, 0, 1), nrow = 2))
91     Y[(nn1 + 1):(nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(-15, -15), matrix(c(1, 0, 0, 1), nrow =
92       2))
93     Y[(nn1 + nn2 + 1):(nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(-15, 15), matrix(c(1, 0, 0,
94       1), nrow = 2))
95     Y[(nh + 1):n, ] <- matrix(runif(p * (n - nh), cc1, 20), ncol = p)
96     if (p > 2) {
97       Y[1:n, 3:p] <- MASS::mvrnorm(n, rep(0, p - 2), diag(p - 2))
98     }
99   } else {
100    Y[1:nn1, ] <- MASS::mvrnorm(nn1, c(5, 5), matrix(c(1, 0, 0, 1), nrow = 2))
101    Y[(nn1 + 1):(nn1 + nn2), ] <- MASS::mvrnorm(nn2, c(-15, -15), matrix(c(1, 0, 0, 1), nrow =
102      2))
103    Y[(nn1 + nn2 + 1):(nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(-15, 15), matrix(c(1, 0, 0,
104      1), nrow = 2))
105    Y[(nh + nn1 + nn2 + 1):(nh + nn1 + nn2 + nn3), ] <- MASS::mvrnorm(nn3, c(15, -15),
106      matrix(c(1, 0, 0, 1), nrow = 2))
107    Y[(2 * nh + 1):n, ] <- matrix(runif((n - 2 * nh) * p, cc1, 20), ncol = p)
108    if (p > 2) {
109      Y[1:n, 3:p] <- MASS::mvrnorm(n, rep(0, p - 2), diag(p - 2))
110    }
111  }
112  true <- rep(0, n)
113  if (G == 3) {
114    true[1:nh] <- rep(1:3, c(nn1, nn2, nn3))
115  } else {
116    true[1:(2 * nh)] <- rep(1:6, c(nn1, nn2, nn3, nn1, nn2, nn3))
117  }
118  return(list(x = Y, true = true))
119 }
```

## 2 Additional simulation results

We have considered 100 random realizations of `simula_tclust()` for the parameter combinations discussed in the manuscript, but with configurations of `balance` and `symmetry` that were not included there. The results obtained in all cases are in perfect agreement with those reported in the manuscript.  $s=1$ ,  $s=2$  and  $s=3$  refer to the 3 scenarios (`scenario=1, 2 and 3`).



Figure A.1: Target function values when `balance=1` and `symmetry=2`.

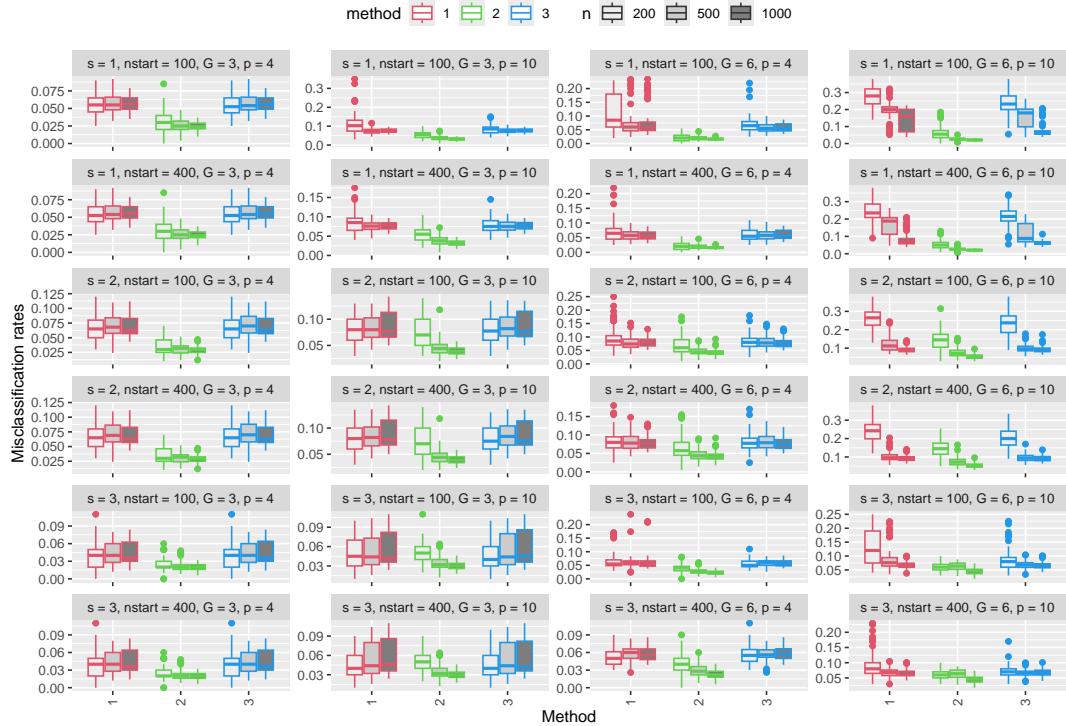


Figure A.2: Misclassification rates when `balance=1` and `symmetry=2`.

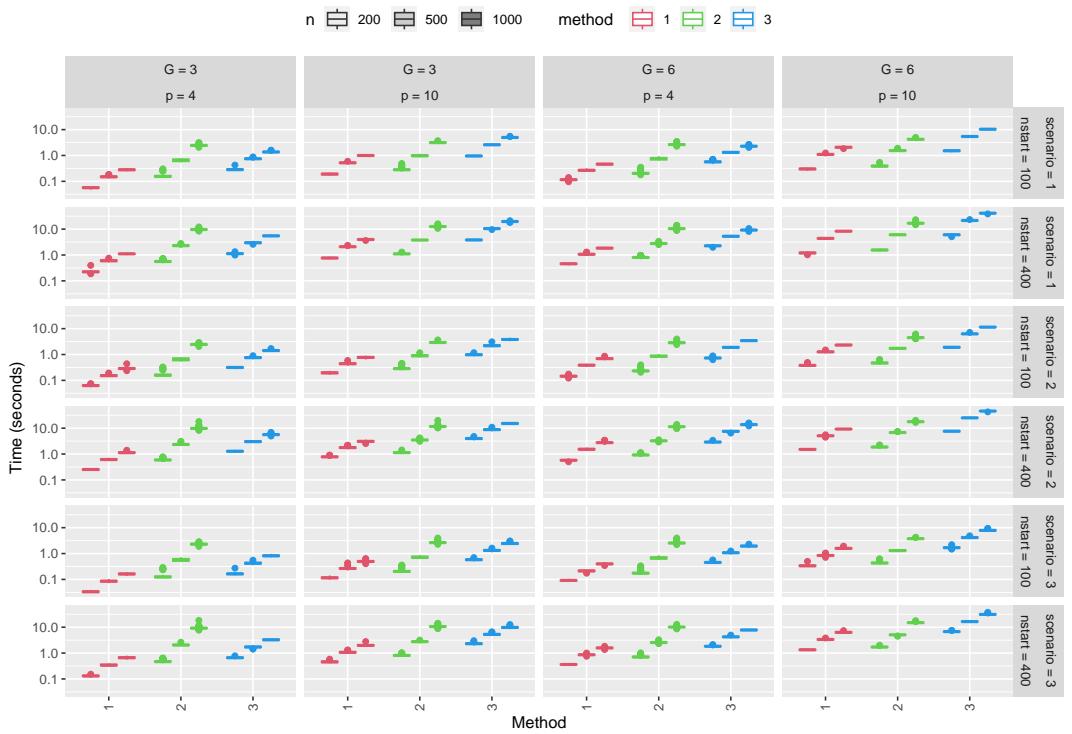


Figure A.3: Computing times in seconds (in logarithmical scale) when **balance**=1 and **symmetry**=2.

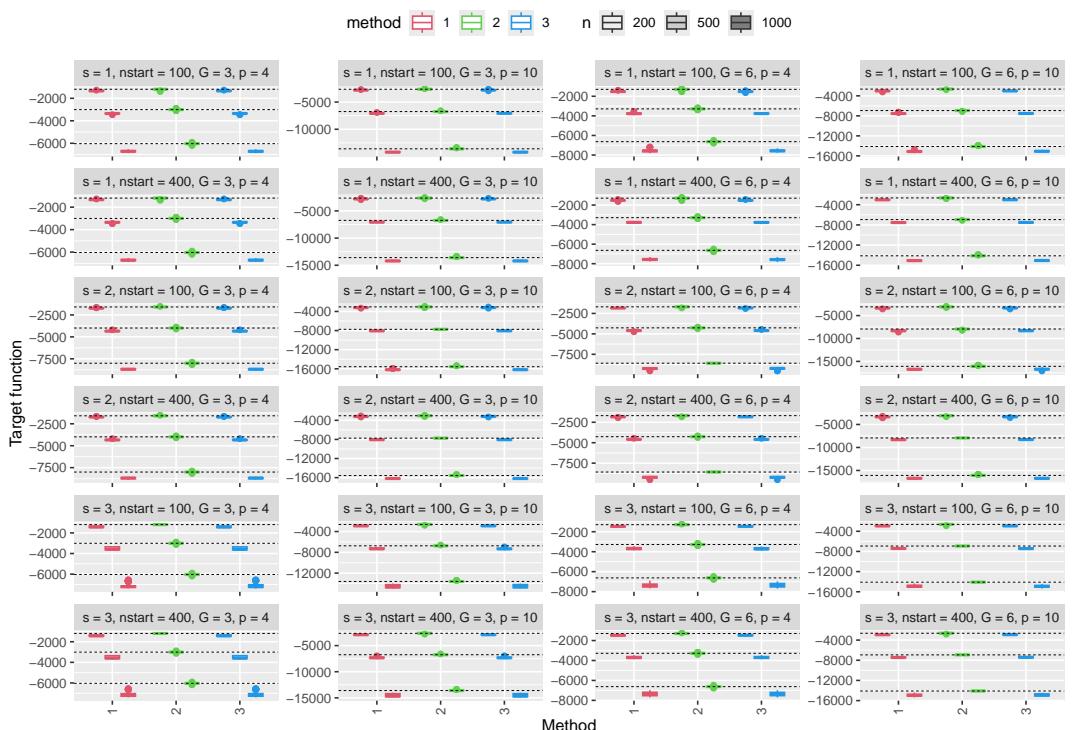


Figure A.4: Target function values when **balance**=2 and **symmetry**=1.

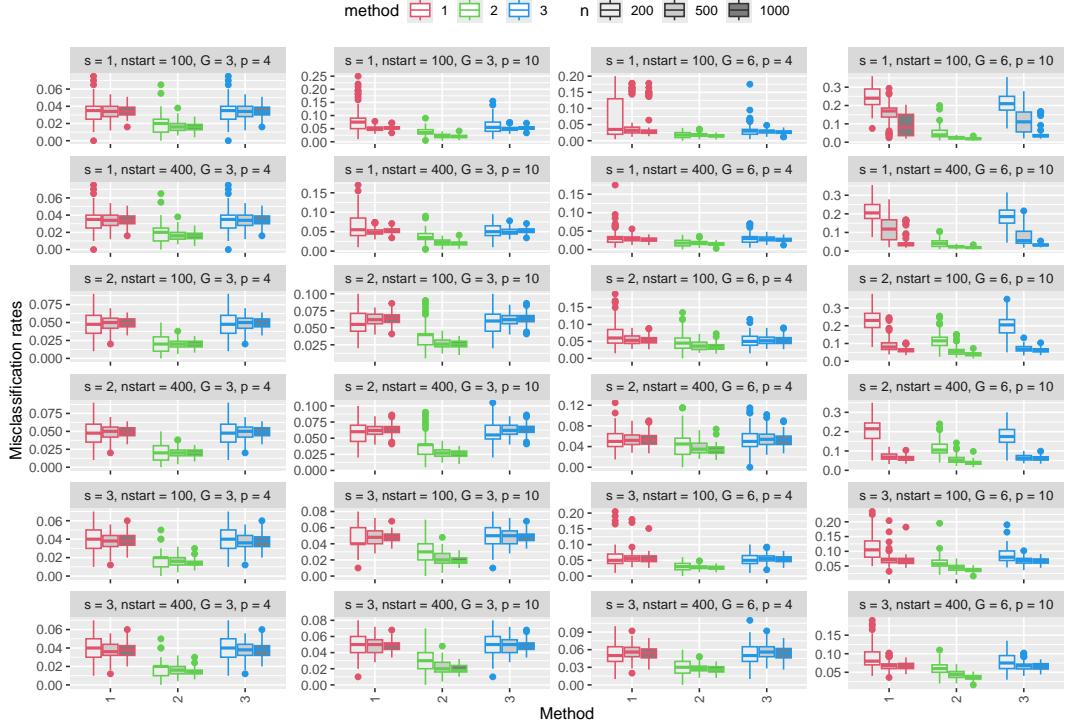


Figure A.5: Misclassification rates when `balance=2` and `symmetry=1`.

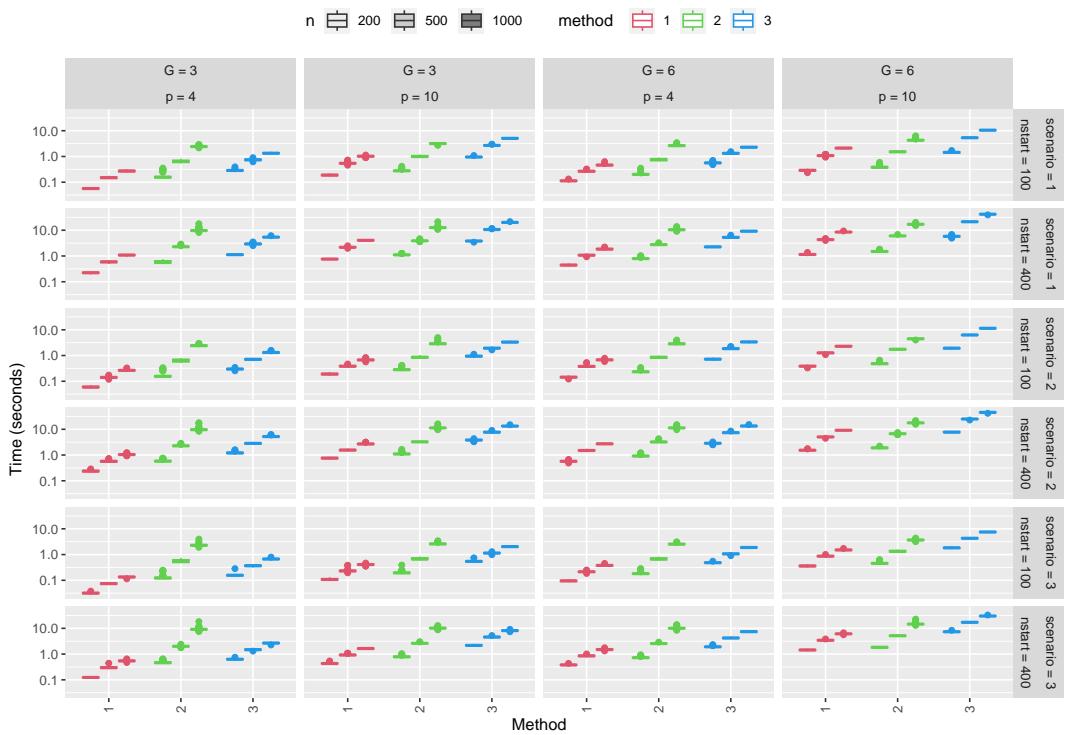


Figure A.6: Computing times in seconds (in logarithmical scale) when `balance=2` and `symmetry=1`.



Figure A.7: Target function values when  $\text{balance}=2$  and  $\text{symmetry}=1$ .

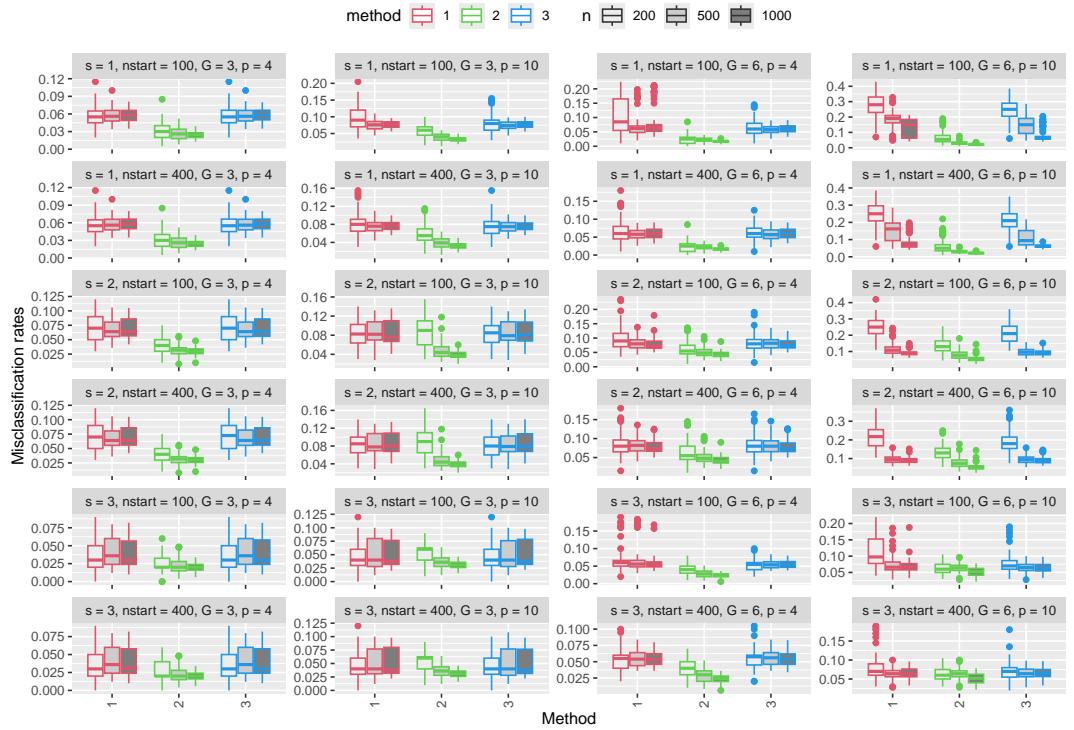


Figure A.8: Misclassification rates when  $\text{balance}=2$  and  $\text{symmetry}=2$ .

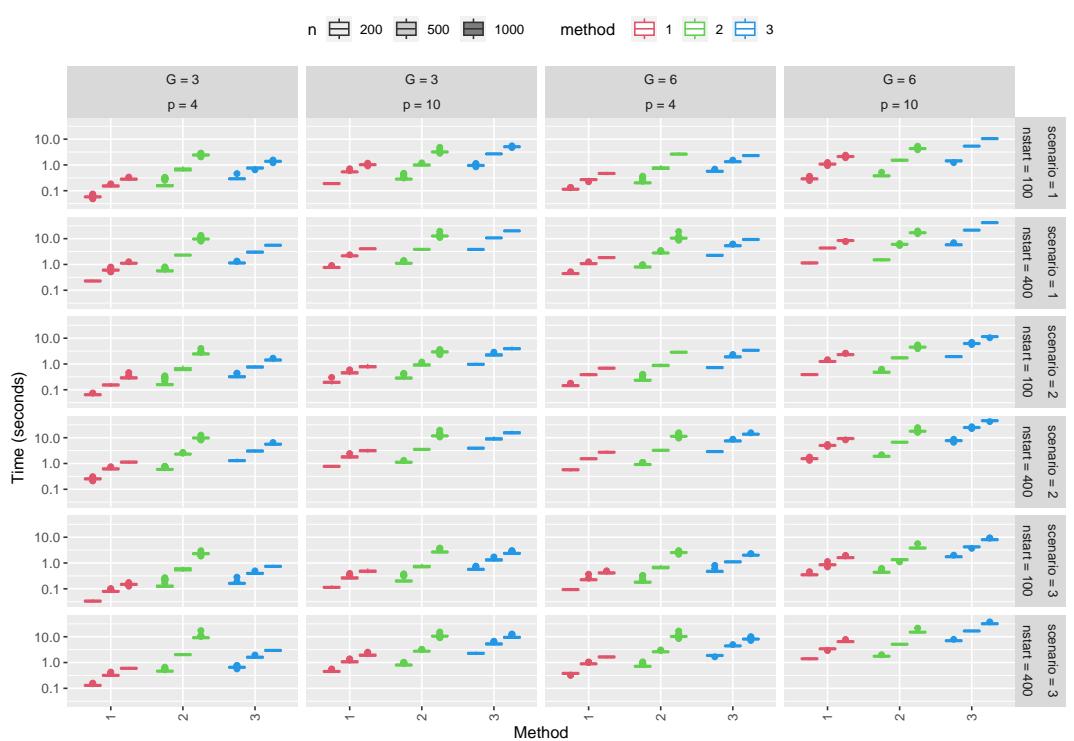


Figure A.9: Computing times in seconds (in logarithmical scale) when `balance`=2 and `symmetry`=2.

### 3 Subsampling and ensemble initialization with $n = 200000$

Results of *Method 1* and *Method 2* for 100 datasets generated by `simula_tclust(n=200000, p, G, scenario, balanced=1, symmetry=1)`. *Method 2* refers to the proposed “Combining subsampling and ensemble initialization” approach with  $n_0 = 800$  (notably smaller than  $n = 200000$ ). Results for *Method 3* with  $5 \times n_{start}$  are not provided due to the substantial computing times required.

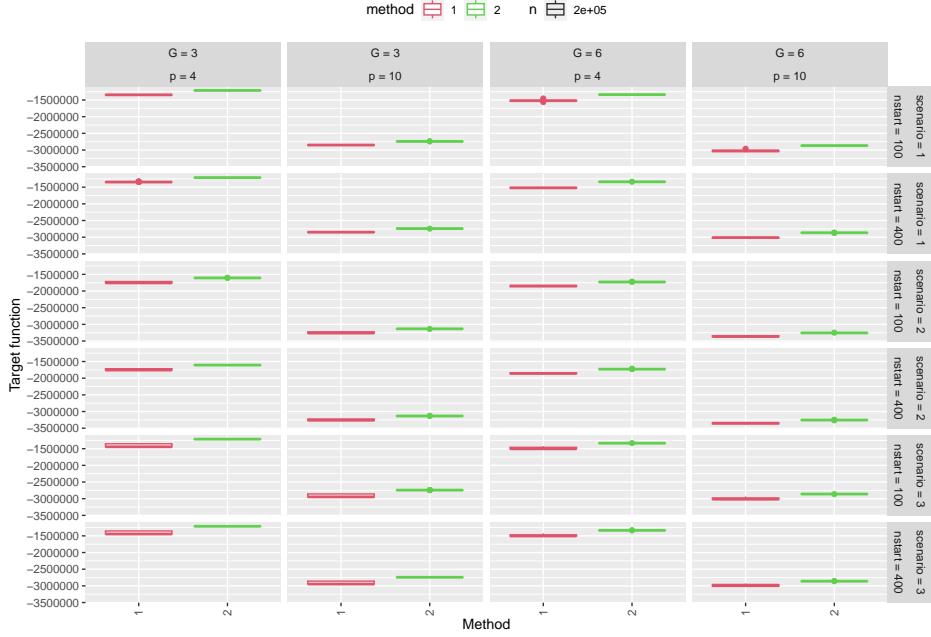


Figure A.10: Target function values when combining subsampling and ensemble initialization with  $n = 200000$ .

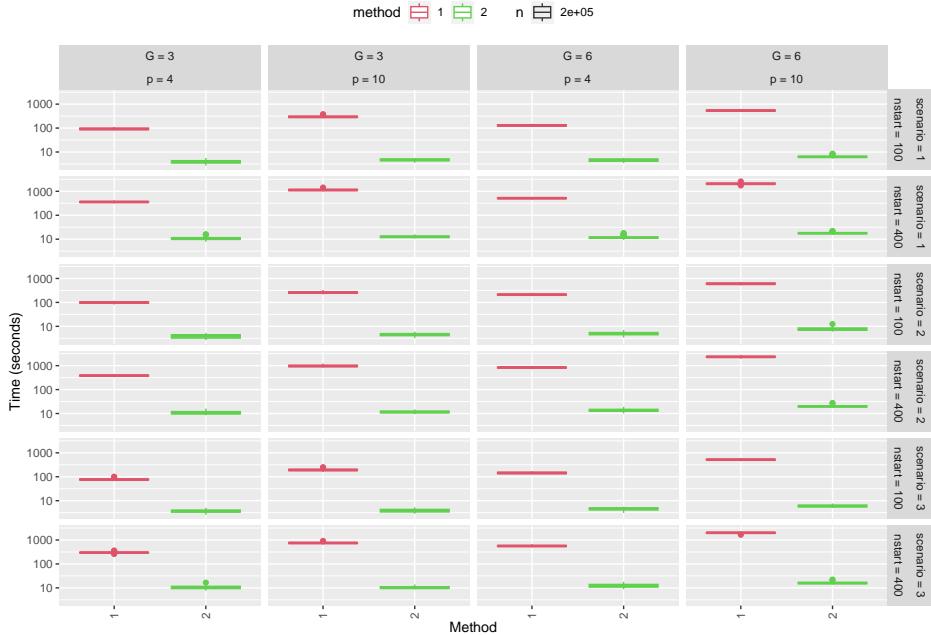


Figure A.11: Computing times (in logarithmic scale) when combining subsampling and ensemble initialization with  $n = 200000$ .

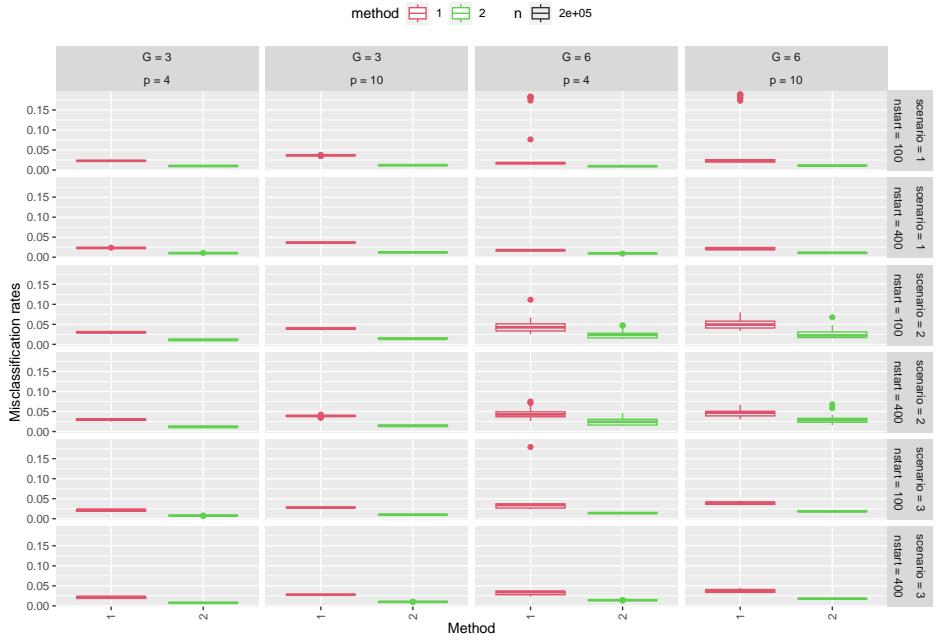


Figure A.12: Misclassification rates when combining subsampling and ensemble initialization with  $n = 200000$ .

## 4 Additional results for the “olive oil” dataset

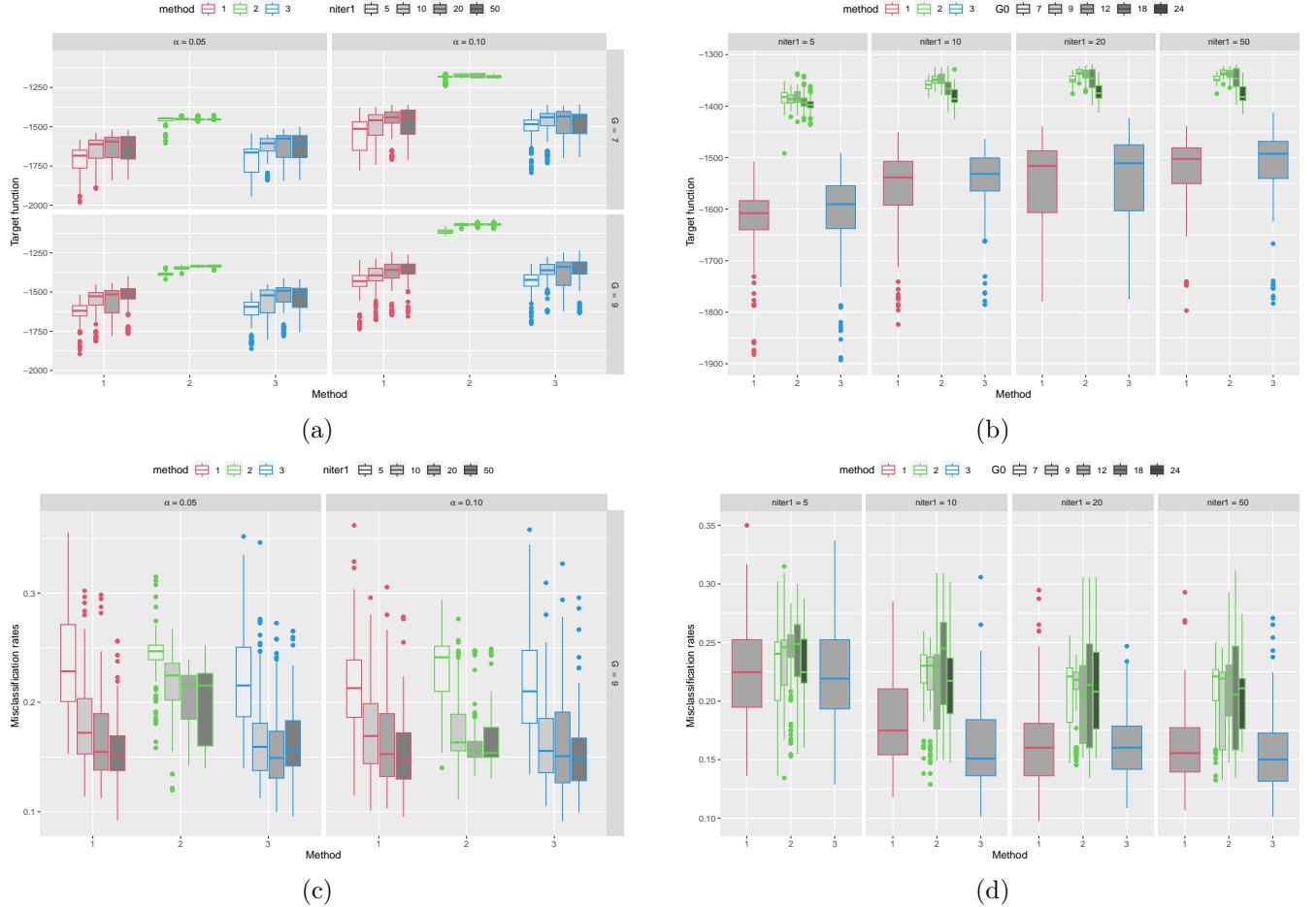


Figure A.13: Results of the experiments as in Section 6 of the manuscript when the restriction factor is set to  $c = 128$ .

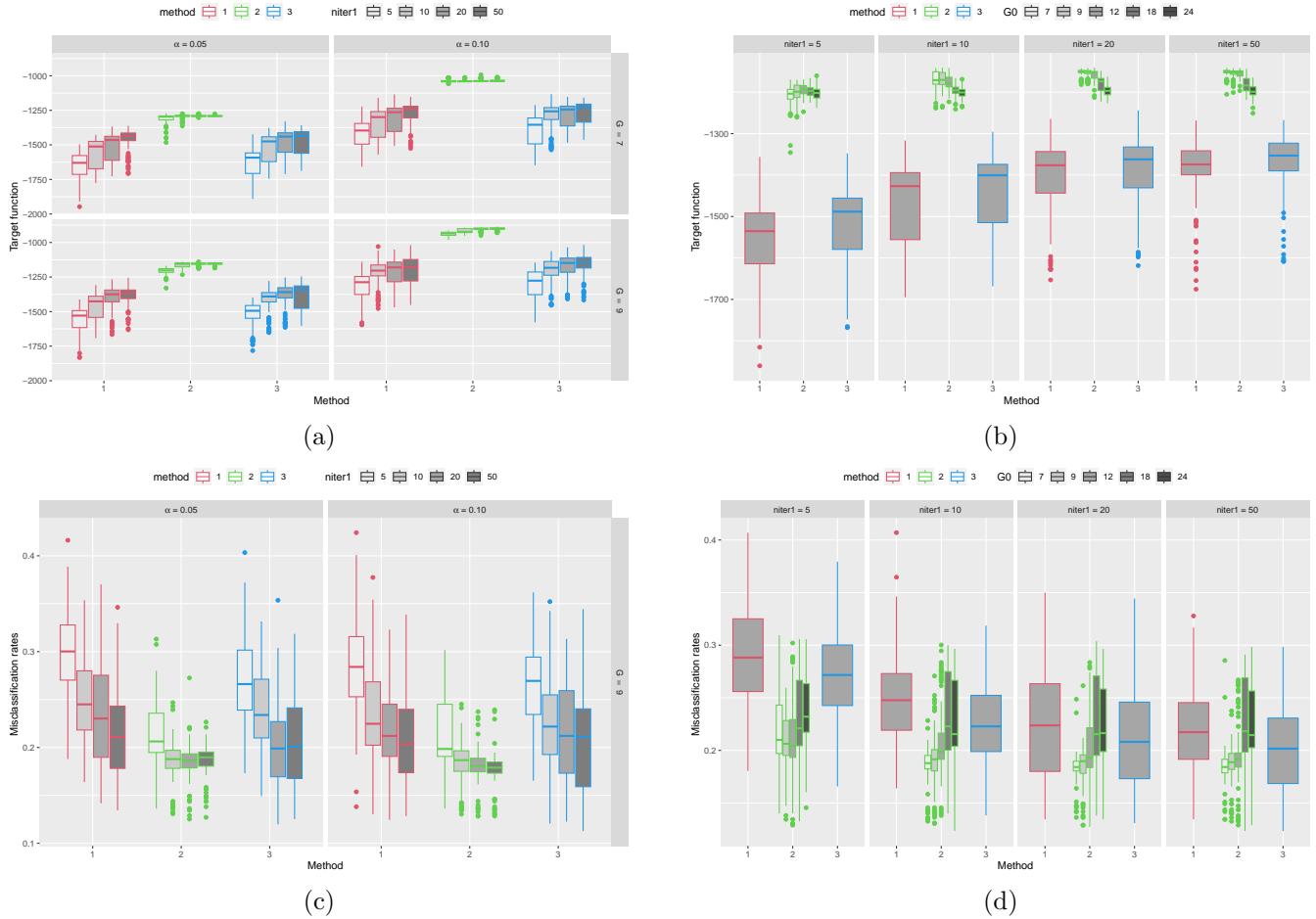


Figure A.14: Results of the experiments as in Section 6 of the manuscript when the restriction factor is set to  $c = 1000$ .