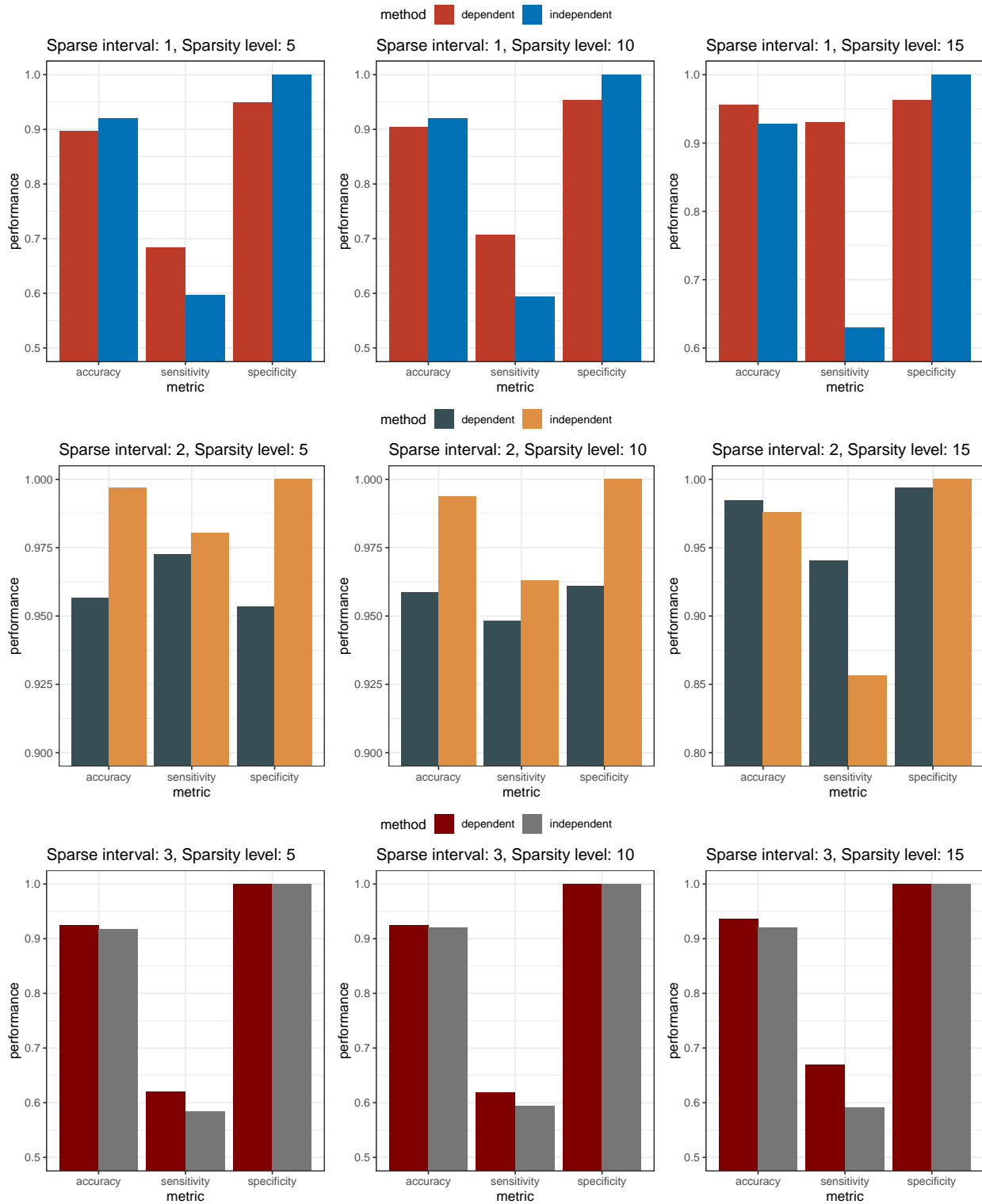


# covdepGE versus competitor

## Summary

I repeated this experiment 9 times. In each experiment, I chose one of the three covariate intervals to be sparse, and the other two to be non-sparse. The non-sparse intervals had 60 individuals in them, while the sparsity level of the sparse interval varied from 5, 10, or 15 individuals. In each experiment, I recorded the following metrics for both the covariate independent and covariate dependent graphical estimation methods:

- Sensitivity: the percentage of edges recovered by the method
- Specificity: the percentage of non-edges recovered by the method
- Accuracy



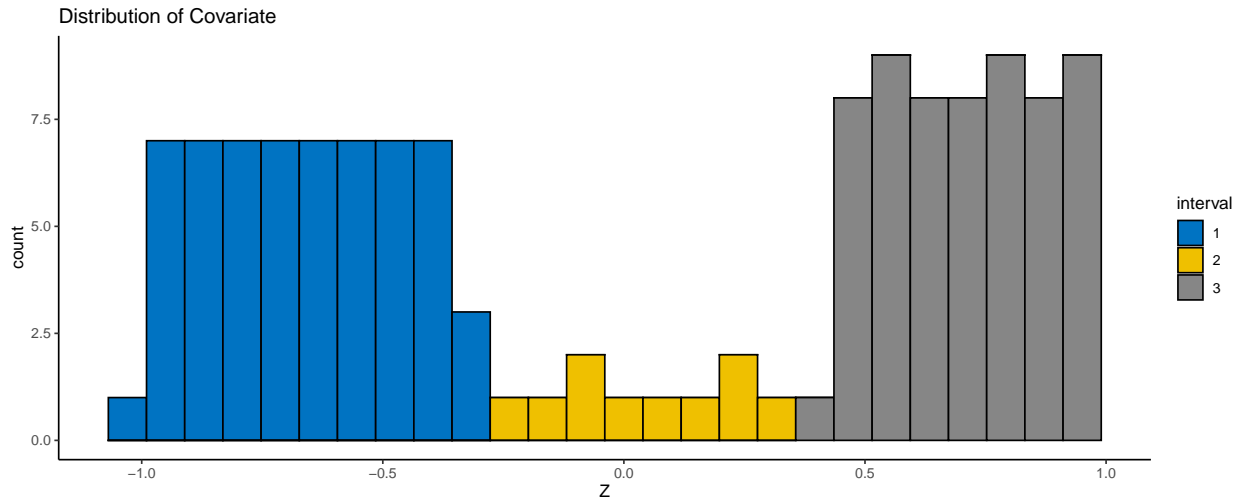
method	sparsity level	sparse interval	sensitivity	specificity	accuracy	unique graphs
dependent	5	1	0.684	0.949	0.896	3
independent	5	1	0.597	1.000	0.920	3
dependent	10	1	0.706	0.953	0.905	3
independent	10	1	0.594	1.000	0.920	3
dependent	15	1	0.930	0.962	0.956	4
independent	15	1	0.630	1.000	0.928	3
dependent	5	2	0.973	0.953	0.956	3
independent	5	2	0.980	1.000	0.997	2
dependent	10	2	0.948	0.961	0.959	4
independent	10	2	0.963	1.000	0.994	2
dependent	15	2	0.940	0.994	0.985	4
independent	15	2	0.856	1.000	0.976	3
dependent	5	3	0.619	1.000	0.924	2
independent	5	3	0.584	1.000	0.917	3
dependent	10	3	0.619	1.000	0.925	2
independent	10	3	0.594	1.000	0.920	3
dependent	15	3	0.670	1.000	0.935	3
independent	15	3	0.591	1.000	0.920	3

Interval	Individual Indices
1	1, ..., 60
2	61, ..., 70
3	71, ..., 130

## Data generation

The extraneous covariate is created as the union of three disjoint intervals with nearly adjacent endpoints. Within each interval, the individuals' covariate values are equally spaced.

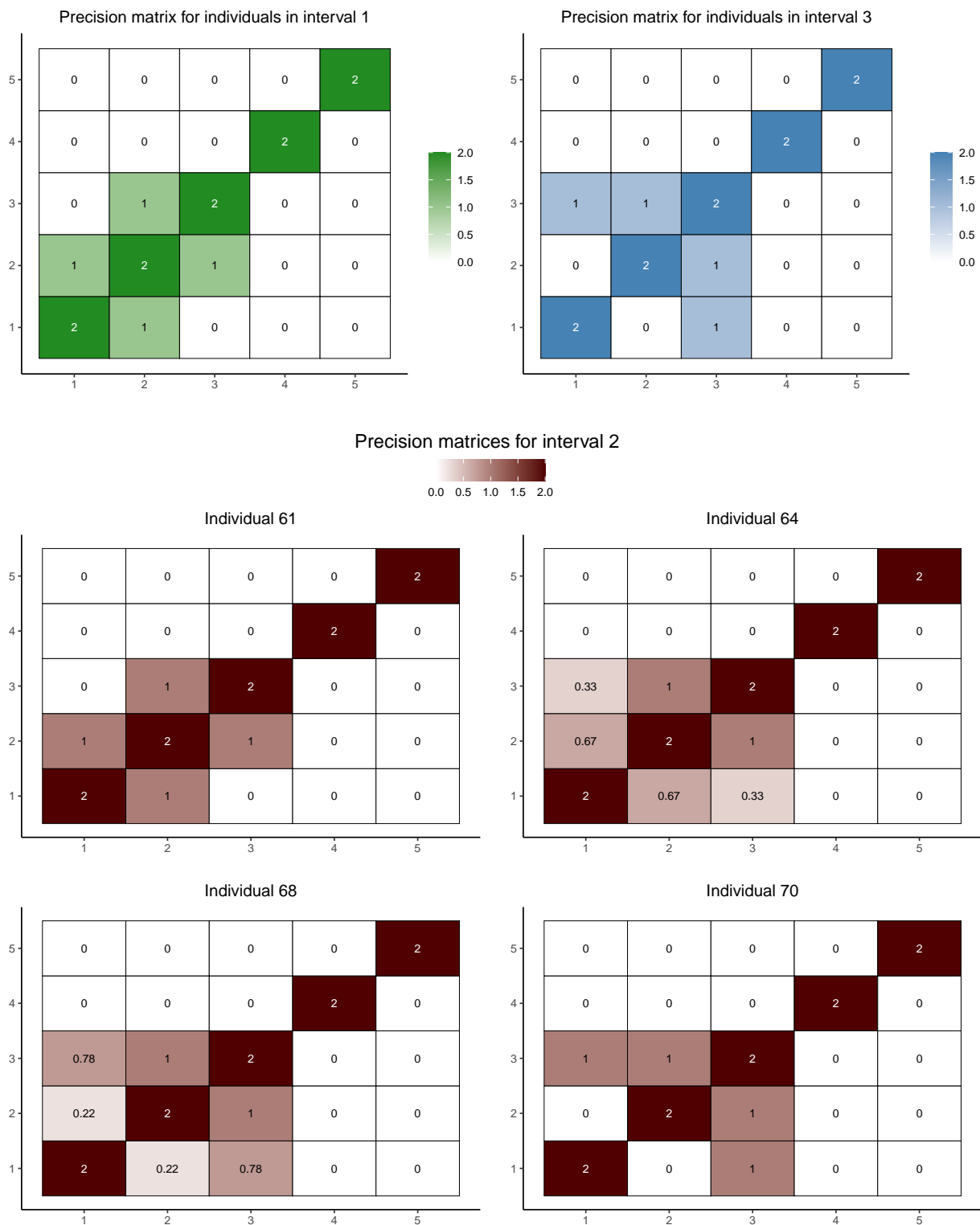
For this demonstration, the second interval is the sparse part of the covariate space, with sparsity level 10. However, in the other 8 experiments, I varied the interval in which the sparsity occurred and the degree of the sparsity.



All of the individuals in interval 1 have the same precision matrix, as do all of the individuals in interval 3.

The first individual in interval 2 has the same precision matrix as those in interval 1.

As the individual index in interval 2 increases, the precision matrix continuously shifts from the precision matrix in interval 1 to the precision matrix in interval 3 such that the last individual in interval 2 has the same precision matrix as the individuals in interval 3.



After creating the precision matrix for each individual, I inverted the matrices to obtain the covariance matrices. I then used the covariance matrices to generate each observation from a 5 dimensional Gaussian distribution centered at  $\vec{0}$ .

## Covariate dependent graph estimation

```
# use varbvs to get the hyperparameter sigma
sigmasq <- sapply(1:(p + 1), function(col_ind) mean(varbvs::varbvs(
  data_mat[, -col_ind], Z, data_mat[, col_ind], verbose = F)$sigma))
sigmasq

## [1] 0.7417301 0.6413059 0.5728140 0.4096810 0.4706576

mean(sigmasq)

## [1] 0.5672377

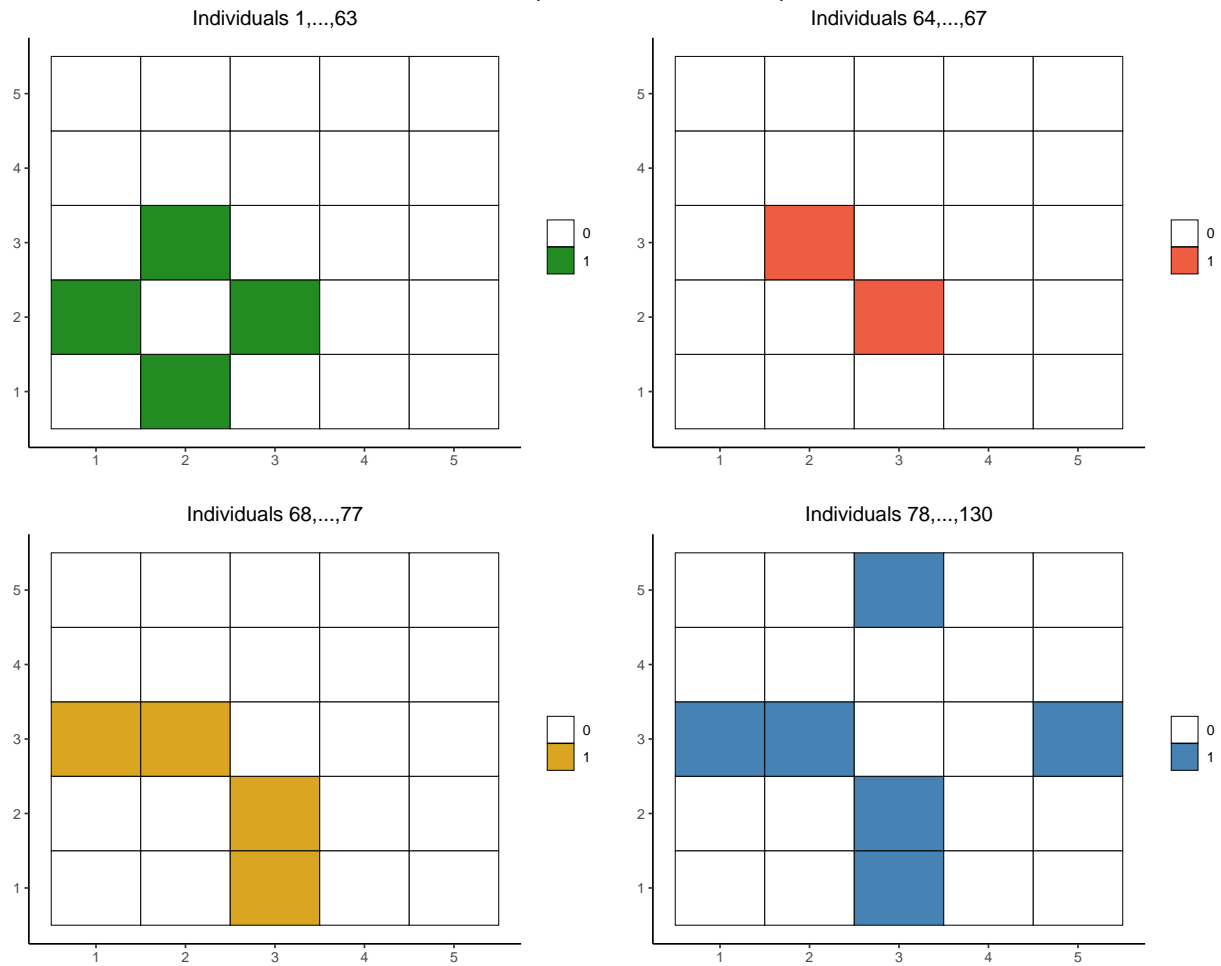
# estimate the covariance structure dependent of the covariate
out_dep <- covdepGE(data_mat,
  Z, # extraneous covariates
  sigmasq = mean(sigmasq), # hyperparameter residual variance
  var_min = 1e-3, # smallest sigmabeta_sq grid value
  var_max = 5, # largest sigmabeta_sq grid value
  n_sigma = 50, # length of the sigmabeta_sq grid
  pi_vec = 0.1, # prior inclusion probability
  tolerance = 1e-10, # variational parameter exit condition 1
  max_iter = 1e3, # variational parameter exit condition 2
  print_time = T
)

## Warning in covdepGE(data_mat, Z, sigmasq = mean(sigmasq), var_min = 0.001, :
## Response 3: 1/50 candidate models did not converge in 1000 iterations

## Warning in covdepGE(data_mat, Z, sigmasq = mean(sigmasq), var_min = 0.001, : For
## 1/5 responses, the selected value of sigmabeta_sq was on the grid boundary. See
## return value VB_details

## Time difference of 23.54996 secs
```

## Covariate Dependent Estimated Graphs



## Covariate independent graph estimation

I first applied Gaussian Mixture Model clustering to the extraneous covariate. The number of clusters is automatically selected by optimizing BIC.

For all of the individuals within each of the clusters identified by GMM, I estimated the shared graph by applying `covdepGE` using a constant value for the extraneous covariate, which will result in the same estimate for all individuals within each cluster.

```
# estimate the dependence structure independent of the covariate

# apply Gaussian Mixture model clustering; selects number of clusters based on
# the model that results in the best BIC
gmm <- Mclust(Z)

# find accuracy of the clustering
fossil::rand.index(gmm$classification, as.numeric(cov_df$interval))
```

```
## [1] 0.9254621
```

```

# find number of clusters in final clustering
(num_clusters <- length(unique(gmm$classification)))

## [1] 2

out_indep <- vector("list", num_clusters)

# iterate over each of the clusters identified by GMM
for (k in 1:num_clusters) {

  # fix the datapoints in the k-th cluster
  data_mat_k <- data_mat[gmm$classification == k, ]

  # use varbvs to get the hyperparameter sigma
  sigmasq_k <- sapply(1:(p + 1), function(col_ind) mean(varbvs::varbvs(
    data_mat_k[, -col_ind], NULL, data_mat_k[, col_ind], verbose = F)$sigma))

  # apply the GGM using covdepGE with constant Z, save the resulting graph
  out_indep[[k]] <- covdepGE(data_mat_k,
    rep(0, nrow(data_mat_k)), # extraneous covariates
    sigmasq = mean(sigmasq_k), # hyperparameter residual variance
    var_min = 1e-3, # smallest sigmabeta_sq grid value
    var_max = 5, # largest sigmabeta_sq grid value
    n_sigma = 50, # length of the sigmabeta_sq grid
    pi_vec = 0.1, # prior inclusion probability
    tolerance = 1e-10, # variational parameter exit condition 1
    max_iter = 1e3, # variational parameter exit condition 2
    print_time = T,
    kde = F, # whether to use kde to calculate bandwidths
    scale = F # whether to scale the extraneous covariates
  )
}

## Warning in covdepGE(data_mat_k, rep(0, nrow(data_mat_k)), sigmasq =
## mean(sigmasq_k), : For 2/5 responses, the selected value of sigmabeta_sq was on
## the grid boundary. See return value VB_details

## Time difference of 3.036036 secs

## Warning in covdepGE(data_mat_k, rep(0, nrow(data_mat_k)), sigmasq =
## mean(sigmasq_k), : Response 1: 7/50 candidate models did not converge in 1000
## iterations

## Warning in covdepGE(data_mat_k, rep(0, nrow(data_mat_k)), sigmasq =
## mean(sigmasq_k), : For 2/5 responses, the selected value of sigmabeta_sq was on
## the grid boundary. See return value VB_details

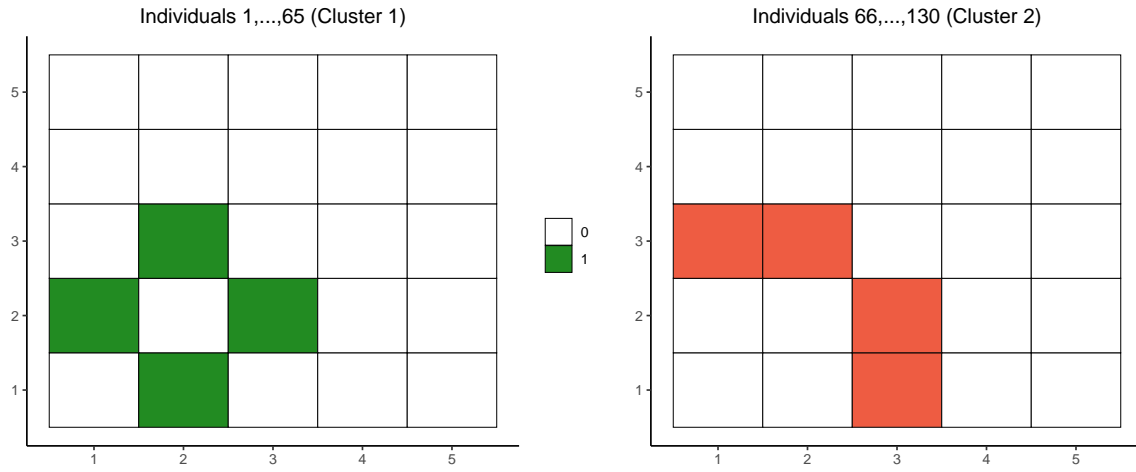
## Time difference of 9.350248 secs

# show the unqiue weights for each of the clusters
lapply(out_indep, function(out) unique(as.vector(out$weights)))

```

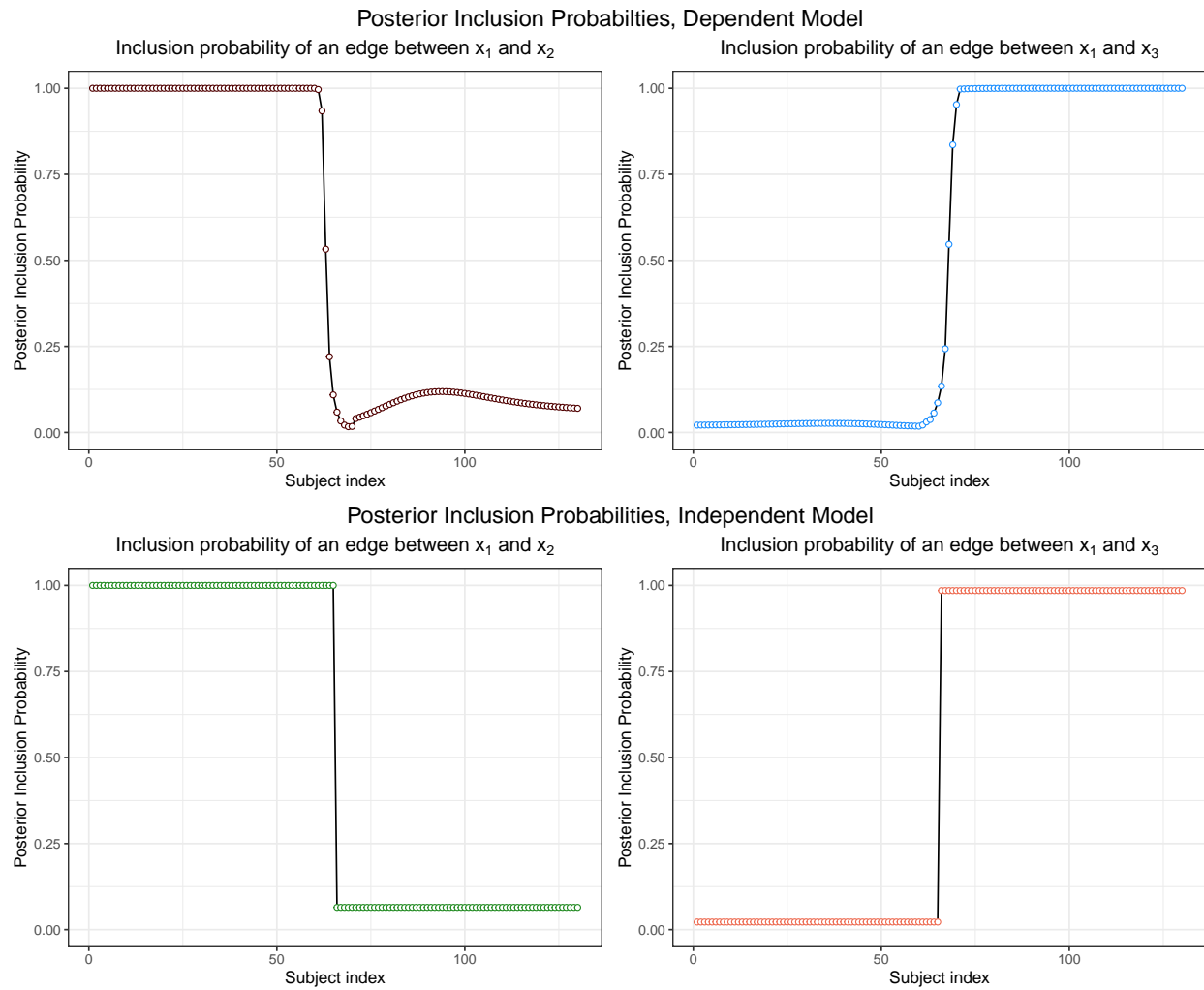
```
## [[1]]
## [1] 1
##
## [[2]]
## [1] 1
```

### Covariate Independent Estimated Graphs

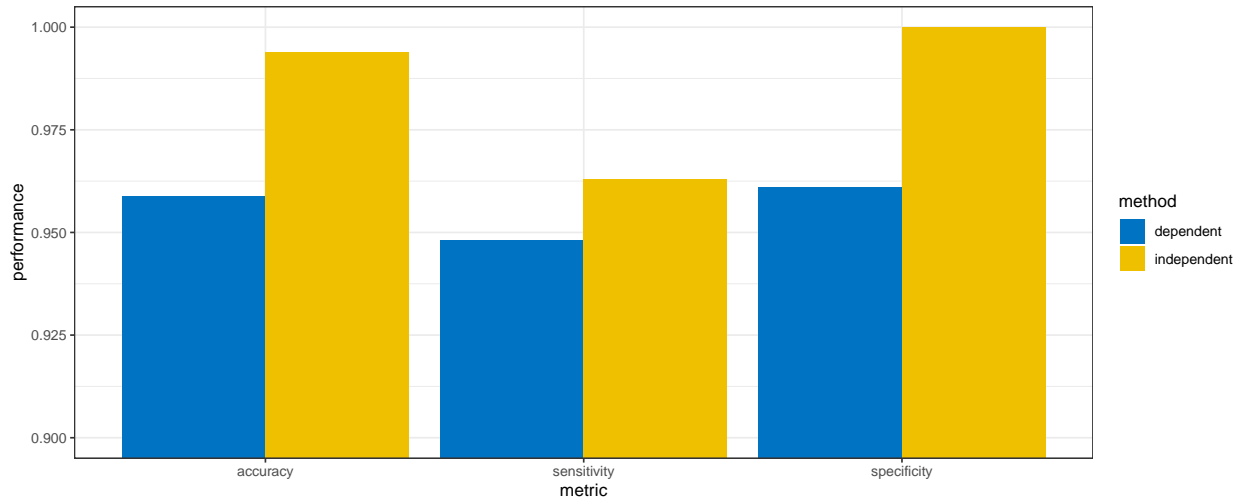




# Performance Analysis



method	metric	performance
dependent	sensitivity	0.948
independent	sensitivity	0.963
dependent	specificity	0.961
independent	specificity	1.000
dependent	accuracy	0.959
independent	accuracy	0.994



## ELBO Monitoring

The latest version of `covdepGE` includes `VB_details` in the return:

`VB_details`: list of  $(p + 1)$  lists; the  $j$ -th list corresponds to the  $j$ -th predictor and contains 6 values:

- `sigma^2_beta`, `pi`: scalars; the final values of `pi` and `sigmabeta_sq` that maximized ELBO over all individuals with the  $j$ -th predictor fixed as the response
- `ELBO`: scalar; the maximum value of ELBO for the final model
- `converged_iter`: scalar; the number of iterations to attain convergence for the final model
- `ELBO_history`: vector; ELBO history by iteration for the final model. If `monitor_final_elbo` is F, then this value will be NULL
- `non_converged`: matrix; each row corresponds to the ELBO history for each of the candidate models that did not converge. If `monitor_cand_elbo` is F, then the ELBO history is omitted, and only the non-convergent `sigmabeta_sq` and `pi` values are provided. If all pairs resulted in convergence, then this value is NULL

```
out <- covdepGE(data_mat, Z, max_iter = 50, tolerance = 1e-10,
  monitor_final_elbo = T, monitor_cand_elbo = T,
  monitor_period = 15, print_time = T)
```

```
## Warning in covdepGE(data_mat, Z, max_iter = 50, tolerance = 1e-10,
## monitor_final_elbo = T, : Response 1: 2/8 candidate models did not converge in
## 50 iterations
```

```
## Warning in covdepGE(data_mat, Z, max_iter = 50, tolerance = 1e-10,
## monitor_final_elbo = T, : Response 2: 4/8 candidate models did not converge in
## 50 iterations
```

```
## Warning in covdepGE(data_mat, Z, max_iter = 50, tolerance = 1e-10,
## monitor_final_elbo = T, : Response 2: final model did not converge in 50
## iterations
```

```
## Warning in covdepGE(data_mat, Z, max_iter = 50, tolerance = 1e-10,
## monitor_final_elbo = T, : For 2/5 responses, the selected value of sigmabeta_sq
## was on the grid boundary. See return value VB_details
```

```
## Time difference of 4.39652 secs
```

```
out$VB_details
```

```
## $'Response 1'
## $'Response 1'$sigma^2_beta'
## [1] 0.5179475
##
## $'Response 1'$pi
## [1] 0.1
##
## $'Response 1'$ELBO
## [1] -9558.61
##
## $'Response 1'$converged_iter
## [1] 34
##
## $'Response 1'$ELBO_history
##           1          16          31          34
## ELBO -11621.12 -9558.611 -9558.61 -9558.61
##
## $'Response 1'$non_converged
##           1          16          31          46          50
## slab var: 0.072, pi: 0.1 -11283.30 -9714.723 -9712.748 -9712.744 -9712.744
## slab var: 0.027, pi: 0.1 -11120.93 -10089.169 -10089.119 -10089.119 -10089.119
##
##
## $'Response 2'
## $'Response 2'$sigma^2_beta'
## [1] 0.5179475
##
## $'Response 2'$pi
## [1] 0.1
##
## $'Response 2'$ELBO
## [1] -10732.39
##
## $'Response 2'$converged_iter
## [1] 50
##
## $'Response 2'$ELBO_history
##           1          16          31          46          50
## ELBO -12787.89 -10732.4 -10732.39 -10732.39 -10732.39
##
## $'Response 2'$non_converged
##           1          16          31          46          50
## slab var: 0.518, pi: 0.1 -12787.89 -10732.40 -10732.39 -10732.39 -10732.39
## slab var: 0.193, pi: 0.1 -12695.03 -10801.36 -10801.23 -10801.23 -10801.23
## slab var: 0.072, pi: 0.1 -12685.18 -11131.58 -11131.50 -11131.50 -11131.50
```

```

## slab var: 0.027, pi: 0.1 -12980.05 -11995.44 -11995.44 -11995.44 -11995.44
##
##
## $'Response 3'
## $'Response 3'$'sigma^2_beta'
## [1] 0.5179475
##
## $'Response 3'$pi
## [1] 0.1
##
## $'Response 3'$ELBO
## [1] -10753.37
##
## $'Response 3'$converged_iter
## [1] 31
##
## $'Response 3'$ELBO_history
##           1           16           31
## ELBO -12987.42 -10753.37 -10753.37
##
## $'Response 3'$non_converged
## NULL
##
##
## $'Response 4'
## $'Response 4'$'sigma^2_beta'
## [1] 0.01
##
## $'Response 4'$pi
## [1] 0.1
##
## $'Response 4'$ELBO
## [1] -6972.793
##
## $'Response 4'$converged_iter
## [1] 8
##
## $'Response 4'$ELBO_history
##           1           8
## ELBO -6972.798 -6972.793
##
## $'Response 4'$non_converged
## NULL
##
##
## $'Response 5'
## $'Response 5'$'sigma^2_beta'
## [1] 0.01
##
## $'Response 5'$pi
## [1] 0.1
##
## $'Response 5'$ELBO
## [1] -8288.941

```

```
##
## $'Response 5'$converged_iter
## [1] 11
##
## $'Response 5'$ELBO_history
##          1          11
## ELBO -8289.062 -8288.941
##
## $'Response 5'$non_converged
## NULL
```