

# Guidelines GORIC(A) benchmarks

Rebecca M. Kuiper

02 April 2024

## Contents

<b>Introduction</b>	<b>1</b>
<b>GORIC(A) weights benchmarks</b>	<b>2</b>
Labelling . . . . .	2
Examples . . . . .	3
<b>Log-likelihood benchmarks</b>	<b>7</b>
Example 3 (ANOVA): Border is true . . . . .	7
Example 1 (ANOVA) Ctd. . . . .	10

## Introduction

To aid qualifying/labeling the height of support for the preferred hypothesis, one can inspect case-specific benchmarks for the GORIC(A) weights (and for the ratio of GORIC(A) weights of the preferred hypothesis versus the other hypotheses in the set). This is especially helpful when you found support for one hypothesis, as opposed to support for the overlap or a boundary/border of two or more hypotheses (cf. ‘Guidelines\_output\_GORIC’ on <https://github.com/rebeccakuiper/Tutorials>). Notably, these benchmarks can also help in noticing support for overlap (if you forgot to check the fit values), since the benchmarks will show that there is maximum/bounded support; as will exemplified in the second example below.

The benchmarks are based on user-specified sets of populations values (as is also used in a power analysis when performing a null hypothesis test). For an ANOVA, the set of populations values can be obtained by specifying population effect size (i.e., Cohen’s  $f$ ) and a specific ratio of population means (which can be the ratio of sample means, that is, the ratio of means found in the data), using the ‘benchmarks\_ANOVA’ function. For all types of statistical models (including ANOVA), one can specify the set of population parameter values themselves, using the ‘benchmarks’ function.

I advise on using the null as a reference population, that is, assuming that one or more equalities (instead of the hypothesized inequalities) hold true in the population. Then, you obtain insight in how (un)likely your results are when one or more equalities are true. The extremer your finding, the more support for your informative, theory-based, inequality-constrained hypothesis.

Additionally, the functions render benchmarks for the ratio of log-likelihood (loglik) weights of the preferred hypothesis versus the other hypotheses in the set and for the (absolute) differences in loglik values of the preferred hypothesis versus the other hypotheses in the set. Note that the log-likelihood benchmarks (under a null population) give insight into the distribution of loglik weight ratios and of the (absolute) loglik differences in case some or all of the group means would be the same. When you calculate these loglik benchmarks for a population in which one or more of the hypothesized inequalities is set to an equality, these loglik benchmarks are helpful in determining support for a boundary hypothesis (if of interest). This will be illustrated in the third example below. Bear in mind that research (e.g., a simulation study) is needed to obtain more insight whether and how well the log-likelihood benchmarks works.

To use the benchmark function, first install and load it:

```
# Install and load benchmark function
if (!require("devtools")) install.packages("devtools")
library(devtools)
install_github("rebeccakuiper/benchmarks")
library(benchmarks)
# For more information regarding the included functions: ?benchmarks
# ?benchmarks_ANOVA
```

## GORIC(A) weights benchmarks

The GORIC(A) weights benchmarks come from several percentiles of sets of GORIC(A) weights assuming that the specified set of populations values is true (for the sample size under consideration). More specifically, the benchmarks are based on the 5%, 35%, 50%, 65%, and 95% percentiles of the GORIC(A) weights for the preferred hypothesis and of the ratios of the GORIC(A) weights for the preferred hypothesis with the other hypotheses.

You can compare your GORIC(A) weight and ratios of GORIC(A) weights of the preferred hypothesis to these benchmarks to make a conclusion regarding the strength of support for the hypothesis (given the assumed set of population parameter estimates and given the sample size). If the benchmarks show a maximum/bounded support (see second example below), there is support for the overlap of two or more hypotheses (which is also signaled by equal log-likelihood values). Otherwise, the GORIC(A) weights (ratios) benchmarks can be used to qualify the support (see first example below).

Additionally, one can check the log-likelihood benchmarks. This can give insight into whether there is support for the boundary of hypotheses (see the third example below, part of the third section).

## Labelling

I am not in favor of cut-off points (or ‘surrounding anchors’), but we need them when we want to label the benchmarks, I propose the following:

benchmark (percentile)	Height support
below 5%	no support
between 5% and 35%	low support
between 35% and 65%	medium support
between 65% and 95%	high (compelling) support
over 95%	very large (tremendous) support

I advise on using some kind of null model as the assumed population. Then, you can see how extreme your finding is (or not) for this null population. The extremer your finding, the more support for your informative, theory-based, inequality-constrained hypothesis.

Alternatively, you may want to use a minimum effect. One option is to specify your hypothesis such that it inspects, for example, minimum differences between parameters (e.g.,  $\mu_1 - \mu_2 > 0.2$  instead of  $\mu_1 > \mu_2$ , that is,  $\mu_1 - \mu_2 > 0$ ). A second option is to use the GORIC(A) weights benchmark for a null model, using a minimum percentile level (for which you think the finding is said to be extreme enough). A third option is to investigate benchmarks using a minimum effect size (or looking at multiple ones). When using the ‘benchmarks\_ANOVA()’ function, you can specify the effect size level(s) (for Cohen’s  $f$ ). When using the ‘benchmarks()’ function, you should specify the population parameter values (reflecting specific effect sizes). Note that the benchmarks differ when you assume different effect sizes or different population parameter estimates (additionally, the benchmarks change with sample size, as do the GORIC(A) weights themselves).

If of interest (e.g., as a sensitivity analysis), you can do this for multiple sets of population parameter estimates; but this may also complicate drawing conclusions (especially when the assumed sets of population parameter estimates differ much).

If your preferred hypothesis does not have the highest fit and you want to inspect benchmarks, I advise on inspecting multiple ratios of population parameter estimates, where some are in agreement with your hypothesis and others in agreement with the data.

Once more, I would inspect populations in which some of the inequalities of your hypothesis/-es of interest are set to equalities (especially if the log-likelihood values seem to be close). This way, you can also check for the support of a boundary hypothesis; as discussed in the third example below.

## Examples

Next, I will discuss two ANOVA examples. More specifically, I will inspect the case-specific benchmarks values (using the ratio of means as in the data). I will look at

- an ANOVA example where we evaluate  $H_1 : \mu_1 > \mu_2 > \mu_3$  versus its complement, and  $H_1$  is true;
- an ANOVA example where we evaluate two overlapping hypotheses, namely  $H_1 : \mu_1 > \mu_2 > \mu_3$  and  $H_2 : \mu_1 > \mu_2, \mu_3$ , together with the unconstrained, and  $H_1$  is true (and thus the others are as well, but they are not the most parsimonious one).

Later on (in another section), I will also discuss the following example:

- an ANOVA example where we evaluate  $H_1 : \mu_1 > \mu_2 > \mu_3$  versus its complement, and the border  $\mu_1 = \mu_2 > \mu_3$  is true. Notably, here, both hypotheses are true,  $H_1$  is the most parsimonious one, and we want to conclude that the border is true.

For a description of interpreting GORIC(A) output, see ‘Guidelines\_output\_GORIC’ on <https://github.com/rebeccakuiper/Tutorials>.

### Example 1 (ANOVA): $H_1$ vs its complement

```
# H1 vs complement - H1 is true
H1 <- "D1 > D2 > D3" # mu1 > mu2 > mu3

# Apply GORIC #
set.seed(123)
results_1c <- goric(fit, hypotheses = list(H1), comparison = "complement")
results_1c
```

restriktor (0.5-50): generalized order-restricted information criterion:

Results:

	model	loglik	penalty	goric	loglik.weights	penalty.weights	goric.weights
1	H1	-155.075	2.833	315.816	0.704	0.697	0.845
2	complement	-155.939	3.667	319.212	0.296	0.303	0.155

---

The order-restricted hypothesis 'H1' has 5.46 times more support than its complement.

```
# Benchmarks based on null pop.es <- c(0) benchmarks_1c <-
# benchmarks_ANOVA(results_1c, pop.es) benchmarks_1c$error.prob.pref.hypo
# benchmarks_1c$benchmarks.weight benchmarks_1c$benchmarks.ratios

# Benchmarks based on null/no, small, medium, and large effect sizes, use:
# pop.es <- c(0, .1, .25, .4) # According to Cohen 1992
pop.es <- c(0, 0.1)
benchmarks_1c <- benchmarks_ANOVA(results_1c, pop.es)
benchmarks_1c$error.prob.pref.hypo
```

```
[1] 0.1547036
```

From the goric output, you can conclude that there is support for  $H_1 : \mu_1 > \mu_2 > \mu_3$ , and  $H_1$  is  $0.85 / 0.15 \approx 5.46$  times more supported than its complement.

The probability that  $H_1$  is not the best is 15.47% (namely, the goric.weight for the complement of  $H_1$ , which is also given by `error.prob.pref.hypo`  $\approx 0.15$ ). This already gives insight into the (un)certainty and, therefore, helps in qualifying the results. Additionally, the benchmarks can help:

Based on the benchmarks, you can check how plausible your finding is (given the assumed population parameter estimates and given the sample size). If you want to compare your results with the situation in which one or more equalities hold true, use a null population (here, an effect size of 0, indicating that the three means are equal). Then, you obtain insight into how (un)likely / how extreme your finding with respect to the inequalities is:

```
benchmarks_1c$benchmarks.weight$`pop.es = 0`
```

	Sample	5%	35%	50%	65%	95%
H1	0.8452964	0.1515391	0.5773783	0.6405995	0.681513	0.7194694

```
benchmarks_1c$benchmarks.ratios$`pop.es = 0`
```

	Sample	5%	35%	50%	65%	95%
H1 vs. H1	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
H1 vs. complement	5.463972	0.1786047	1.366182	1.782412	2.139847	2.564689

When assuming that there is no effect in the population (i.e., `pop.es = 0`), that is, all three group means are equal, the GORIC weight of  $H_1$  (i.e., 0.85) is larger than the 95% percentile. The same holds true for the ratio of GORIC weights of  $H_1$  versus its complement. Hence, our finding is very extreme if the null would be true. Using the table with cut-off values above, this indicates that there is very large (tremendous) support, when assuming no population effect size (given a group size of 40 for each of the 3 groups).

Perhaps you want to inspect a minimum effect. You can either create the hypotheses in such a way that they reflect this, or you can inspect benchmarks for specific effect sizes (or specific population parameter values). Here, we will look at the benchmarks for an effect size (i.e., Cohen's  $f$ ) of 0.1, that is, a small effect size:

```
benchmarks_1c$benchmarks.weight[2]
```

```
$`pop.es = 0.1`
```

	Sample	5%	35%	50%	65%	95%
H1	0.8452964	0.4365551	0.6779699	0.6951248	0.6985532	0.7765229

```
benchmarks_1c$benchmarks.ratios[2]
```

```
$`pop.es = 0.1`
```

	Sample	5%	35%	50%	65%	95%
H1 vs. H1	1.000000	1.000000	1.0000	1.000000	1.000000	1.000000
H1 vs. complement	5.463972	0.7748041	2.1053	2.280031	2.317335	3.474732

The output shows that, when assuming that there is a small effect in the population (i.e., `pop.es = 0.1`) and that the ratio of population means is the same as in the data, the GORIC weight of  $H_1$  (i.e., 0.85) is larger than the 95% percentile (cf. `benchmarks.weight`). The same holds true for the ratio of GORIC weights of  $H_1$  versus its complement (cf. `benchmarks.ratios`). Using the table with cut-off values above, this would indicate that there is very large (tremendous) support, when assuming a small population effect size and a population mean ratio as equal to the sample mean ratio (given a total sample size of  $3 \times 40$ ).

Note that the benchmarks depend on your assumed effect size in the population.

Log-likelihood check:

Before inspecting the height of the support, one may want to establish whether there is support for the overlap / boundary of hypotheses. Since we evaluate an informative hypothesis  $H_1$  versus its complement, we should check whether there is support for a boundary hypothesis (in which one or more inequalities in

$H_1$  is replaced by an equality). For this, one should inspect the log-likelihood / fit values of the hypotheses. When these are close (i.e., the ratio of loglik weights is close to 1 or the absolute difference in loglik values is close to 0), then there is support for (one of their) boundaries. In this case, the loglik values are -155.075 and -156.205, with corresponding loglik.weights of .756 and .244 (and thus a difference of approx. 1.13 and a ratio of approx. 3.10).

Since it is hard to judge what is close, one can inspect the benchmarks for the ratio of log-likelihood (loglik) weights and/or for the (absolute) differences in log-likelihood values. This should then be done for a population in which such a boundary is true. We discuss this in the next section.

For now, we assume that the loglik values are not close.

Population information:

In the data generation, I used a ratio of population means of 3:2:1; implying that  $H_1$  is correct. More specifically, I used population mean values of approximately 0.92, 0.61, and 0.31. This implies that Cohen's  $f$  is .25; thus, there is a medium population effect size (which are in the same order as hypothesized). I then sampled 40 observations for each of the three groups, ran an ANOVA (with three groups), and applied the GORIC. Note that Cohen (1992) suggest that a minimum group size of 52 is needed to find a medium effect when doing null hypothesis testing.

When I would sample more observations, the GORIC(A) weight for  $H_1$  converges to 1 (denoting full support for  $H_1$ ). Note that the benchmarks for the GORIC(A) weight for  $H_1$  will remain the same for a null population and will go to 1 for a non-null population. Note that the error probability then goes to 0 and the the ratio of GORIC(A) weights of  $H_1$  versus its complement then goes to infinity.

## Example 2 (ANOVA): Overlapping hypotheses

```
# H1, H2, and unconstrained (default) - subset/overlap true
H1 <- "D1 > D2 > D3" # mu1 > mu2 > mu3
H2 <- "D1 > D2" # H2: D1 > D2, D3 # mu1 > mu2, mu3

# Apply GORIC #
set.seed(123)
results_12u <- goric(fit, hypotheses = list(H1, H2))
results_12u
```

restriktor (0.5-50): generalized order-restricted information criterion:

Results:

	model	loglik	penalty	goric	loglik.weights	penalty.weights	goric.weights
1	H1	-155.075	2.833	315.816	0.333	0.548	0.548
2	H2	-155.075	3.500	317.149	0.333	0.281	0.281
3	unconstrained	-155.075	4.000	318.149	0.333	0.171	0.171

```
round(results_12u$ratio.gw, 3)
```

	vs. H1	vs. H2	vs. unconstrained
H1	1.000	1.948	3.211
H2	0.513	1.000	1.649
unconstrained	0.311	0.607	1.000

```
# Benchmarks
pop.es <- c(0, 0.1, 0.25)
benchmarks_12u <- benchmarks_ANOVA(results_12u, pop.es)
benchmarks_12u$error.prob.pref.hypo
```

```
[1] 0.1547036
```

```
benchmarks_12u$benchmarks.weight
```

```
$`pop.es = 0`
      Sample      5%      35%      50%      65%      95%
H1 0.5479992 0.1619969 0.4541664 0.5064206 0.5353078 0.5479992
```

```
$`pop.es = 0.1`
      Sample      5%      35%      50%      65%      95%
H1 0.5479992 0.3070648 0.5356692 0.5464097 0.5479992 0.5479992
```

```
$`pop.es = 0.25`
      Sample      5%      35%      50%      65%      95%
H1 0.5479992 0.4489189 0.5479992 0.5479992 0.5479992 0.5479992
```

```
benchmarks_12u$benchmarks.ratios
```

```
$`pop.es = 0`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      1.000000 1.000000 1.000000 1.000000 1.000000 1.000000
H1 vs. H2      1.947734 0.4315439 1.596710 1.888922 1.947734 1.947734
H1 vs. unconstrained 3.211271 0.3490043 1.941405 2.520777 2.941724 3.211271
```

```
$`pop.es = 0.1`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      1.000000 1.000000 1.000000 1.000000 1.000000 1.000000
H1 vs. H2      1.947734 0.7443793 1.935138 1.947734 1.947734 1.947734
H1 vs. unconstrained 3.211271 1.0100021 2.973277 3.183993 3.211271 3.211271
```

```
$`pop.es = 0.25`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      1.000000 1.000000 1.000000 1.000000 1.000000 1.000000
H1 vs. H2      1.947734 1.317257 1.947734 1.947734 1.947734 1.947734
H1 vs. unconstrained 3.211271 2.157690 3.211271 3.211271 3.211271 3.211271
```

From the GORIC output, we can see that both hypotheses are not weak and that  $H_1$  is the best. Additionally, we can see that the log-likelihood values are exactly the same, and thus there is support for the overlap of the hypotheses which is  $H_1$  here (since  $H_1$  is a subset of  $H_2$ ).

If we would now check the GORIC(A) weights benchmarks, we would find again that there is support for the overlap (here,  $H_1$ ) and we should not use the benchmarks: From both types of GORIC(A) weights benchmarks (i.e., `benchmarks.weight` and `benchmarks.ratios`) and when the effect size is higher than 0, we can see that there is maximum support, and our support resembles that. Hence, we conclude that there is support for the overlap. Since there is a maximum support, it is not possible to label this support (notably, one should thus also not compare the results to the benchmarks for a 0 effect size). When you would evaluate the overlap (or the preferred hypothesis) versus its complement, you can inspect the height of the support.

Log-likelihood check:

Here, one clearly finds that the loglik / fit values are exactly the same. Hence, the ratio of loglik weights is exactly 1 and the absolute difference in loglik values is exactly 0. Consequently, there is support for the overlap. In this case, we do not need to check the log-likelihood benchmarks.

Note that the log-likelihood benchmarks (under a null population) give insight into the distribution of loglik weight ratios and of the (absolute) loglik differences in case some or all of the group means would be the same.

Population information:

In the data generation, I used a ratio of population means of 3:2:1; implying that  $H_1$  is correct. More

specifically, I used population mean values of approximately 0.92, 0.61, and 0.31. This implies that Cohen's  $f$  is .25; thus, there is a medium population effect size (which are in the same order as hypothesized). I then sampled 40 observations for each of the three groups, ran an ANOVA (with three groups), and applied the GORIC. Note that Cohen (1992) suggest that a minimum group size of 52 is needed to find a medium effect when doing null hypothesis testing.

When I would sample more observations, it does not (really) affect the GORIC(A) weights for  $H_1$ ,  $H_2$ , and the unconstrained: It will converge to bounds it can take on. The benchmarks for the GORIC(A) weights will also attain the maximum value as will the ratio of weights; and it will for each positive population effect size. For a null population, the GORIC(A) weight benchmarks will remain the same.

## Log-likelihood benchmarks

Before inspecting the height of the support, one may want to establish whether there is support for the overlap / boundary of hypotheses. Since we evaluate an informative hypothesis  $H_1$  versus its complement, we should check whether there is support for a boundary hypothesis (in which one or more inequalities in  $H_1$  is replaced by an equality). For this, one should inspect the log-likelihood / fit values of the hypotheses. When these are close (i.e., the ratio of loglik weights is close to 1 or the absolute difference in loglik values is close to 0), then there is support for (one of their) boundaries.

Since it is hard to judge what is close, one can inspect the benchmarks for the ratio of log-likelihood (loglik) weights and for the (absolute) differences in log-likelihood values. This should then be done for a population in which such a boundary is true.

Next, I will inspect two examples: one in which the border is true (Example 3) and one where it is not (Example 1 continued). I will first use a group size of 40, like in the example above; afterwards, I will inspect a higher sample size to obtain insight into the asymptotic properties of the loglik weights.

For now, I will use that the loglik values are said to be the same if the loglik differences are in between the 5% and 95% percentiles of the loglik benchmarks for a null population. You can of course use a narrower range to be more strict.

Remarks:

- The loglik benchmarks need to be more properly investigated.
- The benchmark function contains four types of loglik benchmarks:
  - \* the loglik ratios (i.e., ratio of loglik weights), which can take on values between 0 and infinity and where 1 means that the loglik values are the same;
  - \* the loglik ratios transformed such that the values lie between 1 and infinity, where 1 means that the loglik values are the same;
  - \* the loglik differences, which can take on values between minus infinity and infinity and where 0 means that the loglik values are the same;
  - \* the absolute loglik differences, which can take on values between 0 and infinity and where 0 means that the loglik values are the same.

Note that the first and the third contain the same information, as do the second and the fourth. Namely, the ratio of loglik weights is the same as the ratio of the likelihood values, and both equal the exponent of the difference in loglik values.

In these guidelines, I will only show the output for the difference between two log-likelihoods.

- One could think about doing a likelihood ratio test (LRT), using a Chi-square distribution or, because of the inequality restrictions, a Chi-bar-square distribution. However, the LRT test is appropriate for nested models, while we are interested in models that are non-nested and even non-overlapping (except for the boundary).

### Example 3 (ANOVA): Border is true

```
# H1 vs complement - border (nl., mu1 = mu2 > mu3) is true
H1 <- "D1 > D2 > D3" # mu1 > mu2 > mu3
```

```
# Apply GORIC #
set.seed(123)
results_1c_border <- goric(fit_border, hypotheses = list(H1), comparison = "complement")
results_1c_border
```

restriktor (0.5-50): generalized order-restricted information criterion:

Results:

	model	loglik	penalty	goric	loglik.weights	penalty.weights	goric.weights
1	H1	-150.549	2.833	306.765	0.482	0.697	0.682
2	complement	-150.477	3.667	308.288	0.518	0.303	0.318

---

The order-restricted hypothesis 'H1' has 2.14 times more support than its complement.

Before we inspect the height of support for the preferred hypotheses, we can check whether there is support for the boundary/border of the two (non-overlapping) hypotheses. By eyeballing, I believe the log-likelihood values are close (-150.549 vs -150.477).

To obtain better evidence for the closeness of the loglik values, I will use the log-likelihood benchmarks functions for several specifications of null populations:

```
## Loglik benchmarks based on null / no effect sizes That is, setting all three
## means equal in the population, using the benchmarks_ANOVA function: pop.es
## <- c(0) benchmarks_1c_border <- benchmarks_ANOVA(results_1c_border, pop.es)
## benchmarks_1c_border$benchmarks.LLratios
## benchmarks_1c_border$benchmarks.LLratios_ge1
## benchmarks_1c_border$benchmarks.difLL
## benchmarks_1c_border$benchmarks.absdifLL

# To obtain more insight into closeness log-likelihood values: Loglik
# benchmarks based on using all possible equalities, using the benchmarks
# function:
est <- coef(fit_border)
pop.est <- matrix(c(mean(est[1:2]), mean(est[1:2]), est[3], est[1], mean(est[2:3]),
  mean(est[2:3]), mean(est[1:3]), mean(est[1:3]), mean(est[1:3])), byrow = TRUE,
  ncol = length(est))
benchmarks_1c_border_allpos <- benchmarks(results_1c_border, pop.est)
# benchmarks_1c_border_allpos$benchmarks.LLratios
# benchmarks_1c_border_allpos$benchmarks.LLratios_ge1
benchmarks_1c_border_allpos$benchmarks.difLL
```

\$`pop.est.nr. = 1`

	Sample	5%	35%	50%	65%	95%
H1 vs. H1	0.00000000	0.000000	0.00000000	0.000000e+00	0.00000000	0.000000
H1 vs. complement	-0.07186354	-1.422054	-0.08842792	-8.319458e-06	0.06165588	1.068567

\$`pop.est.nr. = 2`

	Sample	5%	35%	50%	65%	95%
H1 vs. H1	0.00000000	0.000000	0.00000000	0.0000000000	0.00000000	0.000000
H1 vs. complement	-0.07186354	-1.337118	-0.1225905	-0.009451983	0.01828857	0.63085

\$`pop.est.nr. = 3`

	Sample	5%	35%	50%	65%	95%
H1 vs. H1	0.00000000	0.000000	0.00000000	0.00000000	0.00000000	0.00000000
H1 vs. complement	-0.07186354	-2.419489	-0.4854035	-0.2315981	-0.07077023	0.09057281



```
# benchmarks_1c_border_allpos$benchmarks.absdifLL
```

In this example, the difference in log-likelihood values (-.07) lies between the 35% and the 50% percentile of the loglik benchmarks of the first two null populations; and it is close to the 65% benchmark of the third null population (in which all means are set equal). So, for all three orderings, it is in between the 5% and 95%. Based on this, I conclude that the log-likelihood values of  $H_1$  and its complement are close (i.e., their loglik weight ratio is close to 1 and their (absolute) difference is close to 0). Possibly, one can also conclude that there is a bit more support for a boundary hypothesis in which two means are the same than for one where all three are the same (because the sample value is less extreme for those orderings than for the latter one).

-> ->

Population information:

In the data generation, I used a ratio of population means of 2.5:2.5:1; implying that the boundary of  $H_1$  and its complement is correct (and that  $H_1$  is preferred over its complement, since it is more parsimonious). More specifically, I used population mean values of approximately 0.88, 0.88, and 0.35. This implies that Cohen's  $f$  is .25; thus, there is a medium population effect size. I then sampled 40 observations per group, ran an ANOVA (with three groups), and applied the GORIC.

When I would sample more observations, the GORIC(A) weight for  $H_1$  converges to 1 (denoting full support for  $H_1$ ). Note that the error probability then goes to 0 and the ratio of GORIC(A) weights of  $H_1$  versus its complement then goes to infinity. Nevertheless, this is not of interest now, now we are interested in the closeness of log-likelihood values.

## Higher sample size

Call:

```
lm(formula = y ~ 0 + D, data = sample)
```

Coefficients:

```
      D1      D2      D3
0.8862  0.8609  0.4395
```

```
# Now, group size is 200 (instead of 40)
```

```
# H1 vs complement - border (nl., mu1 = mu2 > mu3) is true
```

```
H1 <- "D1 > D2 > D3" # mu1 > mu2 > mu3
```

```
# Apply GORIC #
```

```
set.seed(123)
```

```
results_1c_border_200 <- goric(fit_border_200, hypotheses = list(H1), comparison = "complement")
```

```
results_1c_border_200
```

restriktor (0.5-50): generalized order-restricted information criterion:

Results:

	model	loglik	penalty	goric	loglik.weights	penalty.weights	goric.weights
1	H1	-829.928	2.833	1665.522	0.509	0.697	0.704
2	complement	-829.962	3.667	1667.258	0.491	0.303	0.296
---							

The order-restricted hypothesis 'H1' has 2.38 times more support than its complement.

We will also now check whether there is support for the border of the two (non-overlapping) hypotheses. By eyeballing, I believe the log-likelihood values are close, but - to obtain more evidence - I will use the log-likelihood benchmarks functions for several specifications of null populations:

```
## Loglik benchmarks based on null / no effect sizes That is, setting all three
## means equal in the population, using the benchmarks_ANOVA function: pop.es
## <- c(0) benchmarks_1c_border_200 <- benchmarks_ANOVA(results_1c_border_200,
## pop.es) benchmarks_1c_border_200$benchmarks.LLratios
## benchmarks_1c_border_200$benchmarks.LLratios_ge1
## benchmarks_1c_border_200$benchmarks.difLL
## benchmarks_1c_border_200$benchmarks.absdifLL

# To obtain more insight into closeness log-likelihood values: Loglik
# benchmarks based on using all possible equalities (and none), using the
# benchmarks function:
est <- coef(fit_border_200)
pop.est <- matrix(c(mean(est[1:2]), mean(est[1:2]), est[3], est[1], mean(est[2:3]),
  mean(est[2:3]), mean(est[1:3]), mean(est[1:3]), mean(est[1:3])), byrow = TRUE,
  ncol = length(est))
benchmarks_1c_border_200_allpos <- benchmarks(results_1c_border_200, pop.est)
# benchmarks_1c_border_200_allpos$benchmarks.LLratios
# benchmarks_1c_border_200_allpos$benchmarks.LLratios_ge1
benchmarks_1c_border_200_allpos$benchmarks.difLL

$`pop.est.nr. = 1`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.0000000 0.000000 0.00000000 0.000000e+00 0.00000000 0.000000
H1 vs. complement 0.03450004 -1.422054 -0.08842792 -8.319458e-06 0.06165588 1.135654

$`pop.est.nr. = 2`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.00000000 0.000000 0.00000000 0.000000e+00 0.00000000 0.00000000
H1 vs. complement 0.03450004 -1.296229 -0.08835953 -5.127627e-05 0.06010068 0.9784858

$`pop.est.nr. = 3`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.00000000 0.000000 0.00000000 0.00000000 0.00000000 0.00000000
H1 vs. complement 0.03450004 -2.419489 -0.4854035 -0.2315981 -0.07077023 0.09057281

# benchmarks_1c_border_200_allpos$benchmarks.absdifLL
```

Also in this example, with a much higher sample size, the difference in log-likelihood values (.04) lies between the 5% and the 95% percentile of the loglik benchmarks for all three orderings. Based on this, I conclude that the log-likelihood values of  $H_1$  and its complement are close (i.e., their loglik weight ratio is close to 1 and their (absolute) difference is close to 0). Additionally, one can argue that there is more support for a boundary hypothesis in which two means are the same than for all three means being the same.

## Example 1 (ANOVA) Ctd.

In this subsection, we look at Example 1 again, where we evaluate  $H_1 \leftarrow "D_1 > D_2 > D_3"$  versus its complement, like in Example 3. In this example,  $H_1$  is true in the population, while in Example 3 the truth is on the border.

In this example, the loglik values are -155.075 and -156.205 with corresponding loglik.weights of .756 and .244. Since it is hard to judge what is close, one can inspect the benchmarks for the ratio of log-likelihood (loglik) weights and for the (absolute) differences in log-likelihood values. This should then be done for a population in which such a boundary is true.

```
# To obtain more insight into closeness log-likelihood values: Benchmarks based
# on using all possible equalities (and none), using the benchmarks function:
est <- coef(fit)
pop.est <- matrix(c(mean(est[1:2]), mean(est[1:2]), est[3], est[1], mean(est[2:3]),
  mean(est[2:3]), mean(est[1:3]), mean(est[1:3]), mean(est[1:3])), byrow = TRUE,
  ncol = length(est))
benchmarks_1c_allpos <- benchmarks(results_1c, pop.est)
# benchmarks_1c_allpos$benchmarks.LLratios
# benchmarks_1c_allpos$benchmarks.LLratios_ge1
benchmarks_1c_allpos$benchmarks.difLL
```

```
$`pop.est.nr. = 1`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.000000 0.000000 0.00000000 0.000000e+00 0.00000000 0.00000000
H1 vs. complement 0.8648427 -1.422054 -0.08842792 -8.319458e-06 0.06001913 0.9639349
```

```
$`pop.est.nr. = 2`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.0000000 0.0000000 0.00000000 0.000000e+00 0.00000000 0.00000000
H1 vs. complement 0.8648427 -1.296229 -0.08355899 4.590777e-05 0.07518101 1.355645
```

```
$`pop.est.nr. = 3`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.0000000 0.0000000 0.00000000 0.00000000 0.00000000 0.00000000
H1 vs. complement 0.8648427 -2.419489 -0.4854035 -0.2315981 -0.07077023 0.09057281
```

```
# benchmarks_1c_allpos$benchmarks.absdifLL
```

In this example, the difference in log-likelihood values (1.13) is larger than the 95% percentile of the loglik benchmarks for all three orderings.

Based on this, I conclude that the log-likelihood values of  $H_1$  and its complement are not close (i.e., their loglik weight ratio is not close to 1 and their (absolute) difference is not close to 0); and thus that there is no support for a boundary hypothesis in which two or three means are the same. Hence, there is no support for a boundary hypothesis, only for  $H_1$ .

## Higher sample size

```
# Now, group size is 200 (instead of 40)
```

```
# H1 vs complement
```

```
H1 <- "D1 > D2 > D3" # mu1 > mu2 > mu3
```

```
# Apply GORIC #
```

```
set.seed(123)
```

```
results_1c_200 <- goric(fit_200, hypotheses = list(H1), comparison = "complement")
```

```
results_1c_200
```

restriktor (0.5-50): generalized order-restricted information criterion:

Results:

	model	loglik	penalty	goric	loglik.weights	penalty.weights	goric.weights
1	H1	-829.928	2.833	1665.522	0.889	0.697	0.949
2	complement	-832.009	3.667	1671.351	0.111	0.303	0.051
---							

The order-restricted hypothesis 'H1' has 18.43 times more support than its complement.

Also here, we will check whether there is support for the border of the two (non-overlapping) hypotheses. By eyeballing, I believe the log-likelihood values are not close. To obtain more evidence, I will use the log-likelihood benchmarks functions for several specifications of null populations:

```
## Loglik benchmarks based on null / no effect sizes That is, setting all three
## means equal in the population, using the benchmarks_ANOVA function: pop.es
## <- c(0) benchmarks_1c_200 <- benchmarks_ANOVA(results_1c_200, pop.es)
## benchmarks_1c_200$benchmarks.LLratios
## benchmarks_1c_200$benchmarks.LLratios_ge1 benchmarks_1c_200$benchmarks.difLL
## benchmarks_1c_200$benchmarks.absdifLL

# To obtain more insight into closeness log-likelihood values: Loglik
# benchmarks based on using all possible equalities (and none), using the
# benchmarks function:
est <- coef(fit_200)
pop.est <- matrix(c(mean(est[1:2]), mean(est[1:2]), est[3], est[1], mean(est[2:3]),
  mean(est[2:3]), mean(est[1:3]), mean(est[1:3]), mean(est[1:3])), byrow = TRUE,
  ncol = length(est))
benchmarks_1c_200_allpos <- benchmarks(results_1c_200, pop.est)
# benchmarks_1c_200_allpos$benchmarks.LLratios
# benchmarks_1c_200_allpos$benchmarks.LLratios_ge1
benchmarks_1c_200_allpos$benchmarks.difLL

$`pop.est.nr. = 1`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.000000 0.000000 0.00000000 0.000000e+00 0.00000000 0.000000
H1 vs. complement 2.080879 -1.422054 -0.08842792 -8.319458e-06 0.06165588 1.104203

$`pop.est.nr. = 2`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.000000 0.000000 0.00000000 0.000000e+00 0.00000000 0.000000
H1 vs. complement 2.080879 -1.296229 -0.08355899 4.590777e-05 0.07518101 1.448344

$`pop.est.nr. = 3`
      Sample      5%      35%      50%      65%      95%
H1 vs. H1      0.000000 0.000000 0.00000000 0.00000000 0.00000000 0.00000000
H1 vs. complement 2.080879 -2.419489 -0.4854035 -0.2315981 -0.07077023 0.09057281

# benchmarks_1c_200_allpos$benchmarks.absdifLL
```

In this example with a higher sample size, the difference in log-likelihood values (2.08) is higher than the 95% percentile of the loglik benchmarks for all three null populations. Based on this, I conclude that the log-likelihood values of  $H_1$  and its complement are not close. Thus, there is no support for a boundary hypothesis, only for  $H_1$ .