

Search for optimal λ 's: a simple but efficient meta-heuristic

Bart JOURQUIN

Université catholique de Louvain, CORE, B-1348 Louvain-la-Neuve, Belgium

bart.jourquin@uclouvain.be

If a systematic test of all the values of λ for a Box-Cox transform of a single independent variable in the range $[-2, 2]$ with a step of 0.1, i.e., 41 different values, is undoubtedly feasible, the exhaustive combinations of λ 's in multivariate cases rapidly becomes a very (too) long computing task because of the combinational nature of the problem. This document proposes a simple, efficient and generic heuristic, useable for a combination of N λ 's. It is tested and validated for cases with $N=1, 2$ and 3.

1 Description

The basic idea of the proposed heuristic is to explore the neighborhood of a given combination of λ 's in N dimensions, starting from an initial “valid” combination (i.e., a solution with the expected sign for all the estimators of the independent variables, all the estimators, including the intercepts, being highly significant (***)).

A specific “hill climbing” algorithm (Russell and Norvig, 2003) is developed. It is a mathematical optimization technique which belongs to the family of local searches. It is an iterative algorithm that starts from an arbitrary solution, and further tries to find better solutions by making incremental changes to the solution until no further improvement can be obtained. Since only convergence to a local maximum can be guaranteed, repeated alternative starting values are usually tested to locate the global maximum, and hence the maximum Log-Likelihood. This is usually referred to as a “shotgun” or “random-restart” hill climbing meta-heuristic (Christian and Griffiths, 2016).

Generally, once all the solutions around an initial solution explored, the strategy used in the hill-climbing algorithms is to pursue the exploration towards the solution that presents the highest improvement (steepest climb). Such an approach implicitly considers that there exists a continuum between all the “valid” solutions. However, this is not always the case for the problem discussed here (see the plots produced by the R scripts provided in this project). Therefore, the algorithm tests all the solutions in each direction and further explores all the solutions with the expected signs having a higher Log-Likelihood than the current solution, even if all the estimators are not highly (***') significant. As multiple search paths are explored, the algorithm must “remember” the solutions to (re)start from. Consequently, the algorithm presented here belongs to the family of hill climbing with backtracking heuristics (Witten et al., 2011).

An important drawback of this strategy is that a same solution has a high probability to be encountered along several search paths. In order to avoid time consuming recomputing of already computed solutions, a hash table - a data structure that implements an associative

array mapping keys (combinations of λ 's) and values (solved model for these λ 's) - is used to store all the already computed results. It is checked each time the result of a logit model is needed during the search. The computing of a λ 's specific logit model is thus only performed when it is not yet present in the hash table.

Beside the definition and the initialization of some global variables (Pseudo-code 1), the heuristic has two major phases: the identification of an initial combination of λ 's to start from and the systematic exploration of its neighborhood until no better solution is found.

<i>N</i>	number of lambdas
<i>range</i>	range to search λ 's in (from -range to +range)
<i>step</i>	interval between 2 successive values of λ
<i>nbDraws</i>	number of valid initial λ combinations to randomly draw
<i>solutionsToExplore</i>	list of solutions to explore around (empty)
<i>bestSolution</i>	best solution found so far (empty)

Pseudo-code 1: Definition of the global variables

To start with (Pseudo-code 2), a set of *nbDraws* combinations of λ 's that produce "valid" solutions are randomly drawn, among which the solution with the highest Log-Likelihood is retained as starting solution for further exploration. Our experience shows that 1 random draw of an initial λ is enough when $N = 1$. For $N = 2$ and $N = 3$ respectively, 5 and 15 draws appear to be sufficient to obtain efficient results (see the "Validation" section below).

```

n = 0
while (n < nbDraws) {
  lambdas = random draw of a combination of N lambdas in range;
  solution = retrieveOrComputeLogit(lambdas);
  if (isValid(solution)) {
    n = n + 1;
    if (solution$logLik > bestSolution$logLik) {
      bestSolution = solution;
    }
  }
}
initialize solutionsToExplore with bestSolution;

```

Pseudo-code 2: Identify a good initial combination of λ '

The exploration around a combination of λ 's can be visualized as in Figure 1. The black dots represent the solution to start with. The other dots are one step backward or forward in each dimension.

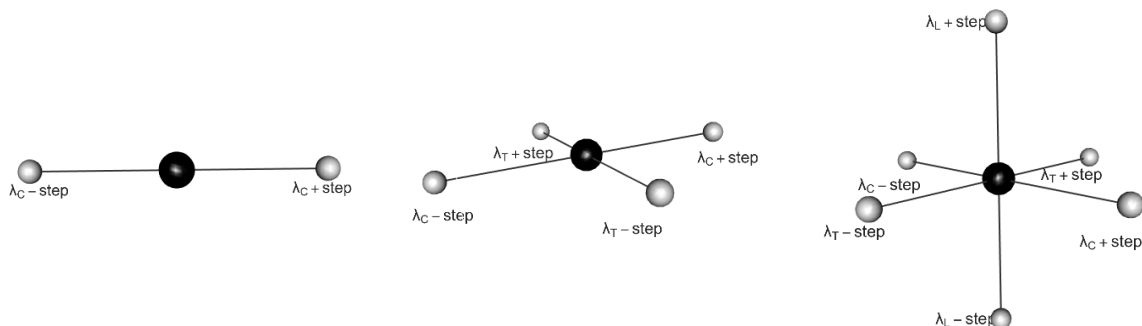


Figure 1: "Explore around" for $N = 1, 2$ and 3

During the exploration process, corresponding to the second phase of the heuristic, every gray dot corresponding to a solution with the expected signs (even if all the estimators of

the solution don't have a '***' signif. code) and having a higher Log-Likelihood than the solution corresponding to the black dot is added to the list of solutions to explore later. All the encountered "valid" solutions are compared to the best one found so far and replace it if better (Pseudo-code 3). Once all the solutions around the black dot are explored, it is removed from the list. Consequently, this list grows and shrinks dynamically, and the exploration ends when the list is empty, meaning that no better "valid" solution can be found (Pseudo-code 4). The particularity of the hill climbing algorithm presented here is that the exploration isn't limited towards the steepest direction, but that all the directions in which the slope increases are tested, as long as the signs of the estimators of the solutions are expected.

Pseudo-code 3 can be implemented as a recursive function (called from itself several times, once for each "dimension" of the combination of λ 's).

```
exploreAround(solution, step) {
  for each lambdas around(step) of solution {
    newSolution = retrieveOrComputeLogit(lambdas);

    if (isValid(newSolution) and newSolution$logLik > bestSolution$logLik) {
      bestSolution = newSolution;
    }

    if (hasExpectedSigns(newSolution) and newSolution$logLik > solution$logLik) {
      add newSolution to solutionsToExplore;
    }
  }
}
```

Pseudo-code 3: Details of the "exploreAround" function

```
repeat {
  solution = first element of solutionsToExplore;
  exploreAround(solution, step);
  remove solution from solutionsToExplore;
} until solutionsToExplore is empty
```

Pseudo-code 4: Explore the neighborhood until no better solution is found

This strategy can further be optimized (Pseudo-code 5) using successive values for *step*, starting from a coarse value and ending with the final granularity of 0.1. Values of 0.4, 0.2 and 0.1 were used for the case presented in this project. This strategy helps to rapidly converge towards the neighborhood of a good solution using large steps, then exploring the surrounding with gradually smaller steps.

```
for each currentStep in (0.4, 0.2, 0.1) {
  repeat {
    solution = first element of solutionsToExplore;
    exploreAround(solution, currentStep);
    remove solution from solutionsToExplore;
  } until solutionsToExplore is empty

  initialize solutionsToExplore with bestSolution;
}
```

Pseudo-code 5: Explore the neighborhood until no better solution is found (optimized version)

The two phases are repeated several times (shotgun meta-heuristic with 3 attempts in the exercise presented in this document) in order to try to locate the global maximum. Again,

the use of the hash table presented before helps to limit the number of logits to compute as, from run to run, already computed solutions can be directly fetched.

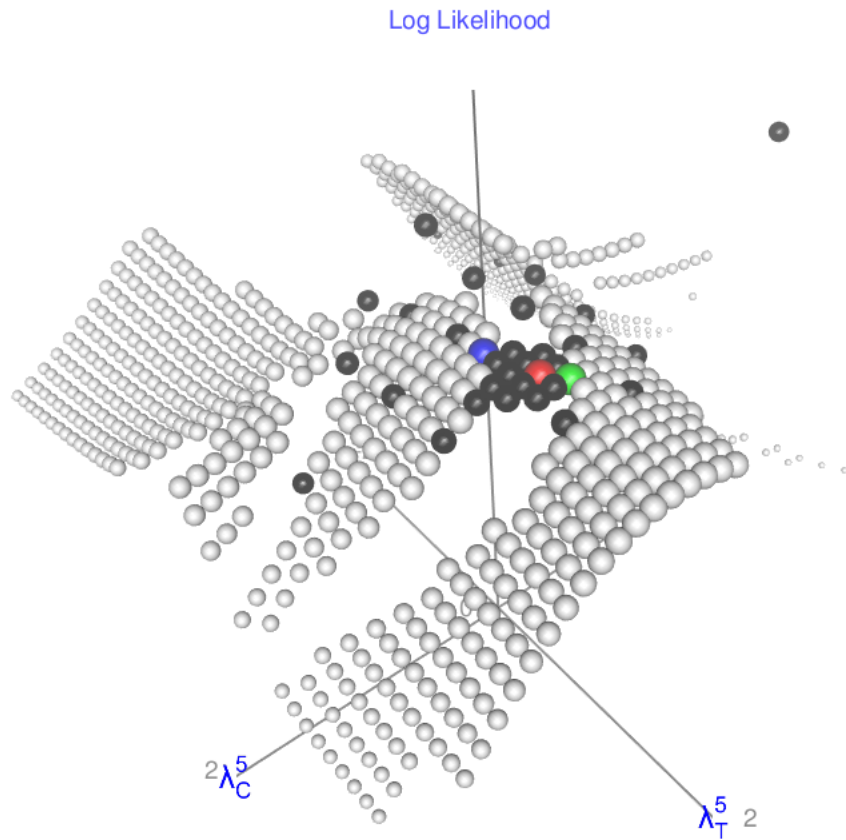


Figure 2: Set of solutions tested by the heuristic (bivariate case, NST-R 5)

Figure 2 gives an idea on how the heuristic finds its path to a solution. The example given here is a bivariate case, because those representing problems with three variables are too cluttered to be presented in this document. This example has been chosen because it doesn't illustrate a "hit" as the solution found by the heuristic (blue dot) doesn't correspond to the exact solution under constraint (green dot). All the black dots represent the solutions tested by the heuristic. Some seem to be in the middle of nowhere; they correspond to random draws of the first phase, that are not "valid" solutions. The density of black dots becomes higher in the neighborhood of the solution. This is a result of the "optimized" approach presented in Pseudo-code 5. In this example, only 34 logit computations were needed (instead of 1 641) to identify a solution (blue dot), which Log-Likelihood is very close to the best solution (green dot). The heuristic didn't converge to the green dot because the red dot (unconstrained max Log-Likelihood) and its surrounding solutions are not "valid" solutions.

1.1 Validation

The results of the heuristic can be compared to the exact solutions computed for the dataset using the brute force algorithm that tests all the possible combinations of λ 's.

As the heuristic starts from a randomly drawn solution, it was run 50 times in order to measure its average and worst-case performances.

The next 3 tables summarize the performance indicators for the univariate, bivariate and trivariate cases. For each group of commodities, one finds:

- The average amount of logit computations needed by the heuristic;
- The largest amount of logit computations (worst case) encountered during the 50 runs;
- The average number of “checks”, i.e., retrievals of already computed solutions;
- The number of “hits” (when the solution of the heuristic corresponds to the exact one);
- The relative difference (in %) between the best and the worst Log-Likelihood among all the “valid” solutions identified by the brute force algorithm. This indicator helps to evaluate the figures of the last two columns;
- The average difference (in %) between the Log-Likelihood of the solution found by the heuristic and the one of the exact solution;
- The largest (worst) difference (in %) between the Log-Likelihood of the solution found by the heuristic and the one of the exact solution.

NST-R	Logit computations (brute force = 41)		Avg checks	Hits (50)	Log-Likelihood (%)		
	Average	Worst			Max Δ (exact)	Avg Δ with best	Worst case Δ with best
0	13	18	20	50	-15%	0%	0%
1	12	16	19	50	-19%	0%	0%
2	13	18	19	50	-18%	0%	0%
3	13	17	20	50	-35%	0%	0%
4	13	16	19	50	-28%	0%	0%
5	13	17	20	50	-8%	0%	0%
6	13	17	20	50	-14%	0%	0%
7	13	17	19	50	-9%	0%	0%
8	10	13	19	50	-15%	0%	0%
9	13	18	19	50	-3%	0%	0%

Table 1: Performances indicators for the univariate case

It comes out (Table 1) that, when only one λ must be identified, the heuristic always converges to the exact solution, but it only needs, on average, 30% of the logit computing's needed by the brute force algorithm. This is interesting but even the brute force computing time is not really a problem here. Given the number of possible λ 's combinations for multivariate cases, it is more interesting to analyze the way in which the heuristic finds its path to a good result in more complex distributions of “valid” solutions.

When the heuristic is applied to the bivariate case (Table 2), it converges after an average of 64 logit computations, i.e., less than 4% of what is needed to obtain the exact solution with a brute force approach. Beside the logit computations, an average 82 retrievals of already computed solutions are also needed. Even if the final solution can differ from the exact one (cfr. the number of hits), its Log-Likelihood is always close to the best one, even for the worst cases.

Finally, Table 3 shows that, when applied to the trivariate case, the heuristic finds a solution after an average of 408 logit computations (and 596 checks), which represents only 0,5% of

the runs needed to find the exact solution with the brute force algorithm. The computing time is thus 200 times faster. The heuristic clearly breaks the combinational logic of the problem. Nevertheless, the Log-Likelihoods of the solutions are very close to the best ones, even for the worst cases.

NST-R	Logit computations (brute force = 1,681)		Avg checks	Hits (50)	Log-Likelihood (%)		
	Average	Worst			Max Δ (exact)	Avg Δ with best	Worst case Δ with best
0	62	83	52	50	-13%	0%	0%
1	58	100	90	34	-14%	0%	0%
2	79	110	65	44	-14%	-0.2%	-2%
3	59	81	57	50	-30%	0%	0%
4	57	80	51	11	-28%	0%	0%
5	70	106	202	14	-7%	0%	0%
6	54	72	46	50	-15%	0%	0%
7	74	124	106	18	-5%	0%	0%
8	65	112	93	49	-12%	0%	-1%
9	66	89	60	50	-5%	0%	0%

Table 2: Performances indicators for the bivariate case

NST-R	Logit computations (brute force = 68,921)		Avg checks	Hits (50)	Log-Likelihood (%)		
	Average	Worst			Max Δ (exact)	Avg Δ with best	Worst case Δ with best
0	320	392	156	46	-12%	0%	0%
1	256	410	688	46	-14%	0%	-1%
2	627	810	1112	4	-12%	-1.0%	-2%
3	606	758	182	50	-28%	0%	0%
4	740	1054	830	1	-26%	-0.9%	-2%
5	249	352	437	50	-6%	0%	0%
6	305	376	355	29	-15%	0%	0%
7	384	560	1217	0	-5%	0%	0%
8	237	320	354	43	-12%	0%	0%
9	359	509	628	28	-5%	-0.1%	-1%

Table 3: Performances indicators for the trivariate case

NST-R	Logit computations (brute force = 68,921)		Avg checks	Hits (50)	Log-Likelihood (%)		
	Average	Worst			Max Δ (exact)	Avg Δ with best	Worst case Δ with best
0	142	224	138	45	-12%	0%	0%
1	132	214	373	47	-14%	0%	0%
2	244	317	569	1	-12%	-1.8%	-5%
3	235	408	136	38	-28%	0%	0%
4	284	424	301	1	-26%	-1.1%	-3%
5	133	240	436	45	-6%	0%	0%
6	138	203	217	7	-15%	0%	-1%
7	192	307	500	0	-5%	0%	-1%
8	109	187	161	36	-12%	-0.1%	-2%
9	160	221	469	20	-5%	-0.3%	-1%

Table 4: Performances indicators for the trivariate case (1 shot heuristic)

The heuristic was also run 50 times with a single shotgun (instead of 3). The results for the trivariate case are published in Table 4 and must be compared to those of Table 3. It appears that the results are very similar, while the average amount of logit computations and

retrievals of already computed solutions are respectively equal to 178 and 330, i.e., roughly 50% of what was needed with 3 shotguns. A single shotgun could thus be considered as enough if computing time is an important constraint.

2 References

Christian, B., Griffiths, T., 2016. Algorithms to live by: The computer science of human decisions, First U.S. Edition. ed. Henry Holt and Company, New York.

Russell, S.J., Norvig, P., 2003. Artificial intelligence: a modern approach, 2nd ed. ed, Prentice Hall series in artificial intelligence. Prentice Hall/Pearson Education, Upper Saddle River, N.J.

Witten, I.H., Frank, E., Hall, M.A., 2011. Data mining: practical machine learning tools and techniques, 3rd ed. ed, Morgan Kaufmann series in data management systems. Morgan Kaufmann, Burlington, MA.