Replicability and Reproducibility in Evolutionary Optimization

Luís Paquete

paquete@dei.uc.pt http://www.uc.pt/go/paquete

University of Coimbra, Portugal



Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, nequires prior specific permission and/or a fee. Request permissions from Permission Badeacn. org. GECCO '21 Companion, July 10-14, 2021, Lille, France @ 2021 Copyright is held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-14503-8351-6/21/07. https://doi.org/10.1145/3448726.3461405



manuel.lopez-ibanez@uma.es http://lopez-ibanez.eu

University of Málaga, Spain





Main reference

Manuel López-Ibáñez, Juergen Branke and Luís Paquete. **Reproducibility in Evolutionary Computation.** *Arxiv preprint arXiv:20102.03380 [cs.Al], 2021.* https://arxiv.org/abs/2102.03380

Instructors

Luís Paquete is Associate Professor at the Department of Informatics Engineering, University of Coimbra, Portugal. He received his Ph.D. in Computer Science from the Technical University of Darmstadt, Germany, in 2005 and a M.S. in Systems Engineering and Computer Science from the University of Algarve, Portugal, in 2001. His research interest is mainly focused on exact and heuristic solution methods for multiobjective combinatorial optimization problems. He is in editorial board of Operations Research Perspectives and Area Editor at ACM Transactions on Evolutionary Learning and Optimization.



http://www.uc.pt/go/paquete

Manuel López-Ibáñez is a "Beatriz Galindo" Senior Distinguished Researcher at University of Málaga since 2020 and a Senior Lecturer (Assistant Professor) at the Alliance Manchester Business School, University of Manchester, UK, since 2015. He received the M.S. degree in computer science from the University of Granada, Granada, Spain, in 2004, and the Ph.D. degree from Edinburgh Napier University, U.K., in 2009. Between 2011 and 2015, he was a Postdoctoral Researcher of the Belgian F.R.S.-FNRS at the IRIDIA laboratory in the Université Libre de Bruxelles (ULB), Brussels, Belgium.



http://lopez-ibanez.eu

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization



Overview

Computer science is in part a scientific discipline concerned with the empirical study of a class of phenomena, in part a mathematical discipline concerned with the formal properties of certain classes of abstract structures, and in part a technological discipline concerned with the cost-effective design and construction of commercially and socially valuable products

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

[Wegner, 1976]

[loannidis, 2005]

EC as an empirical scientific discipline

Scientific Method Observe a phenomenon EAX crossover shows local optimisation in the TSP [Nagata & Kobayashi, 1997] Construct a hypothesis EA + EAX produces better solutions for the TSP Conduct an experiment Draw conclusion about hypothesis: either provisionally accepted or it is *falsified* (with some statistical confidence)

Falsifiability + Reproducibility

 \Rightarrow build research community consensus

⇒ "Laws of qualitative structure" [Newell & Simon, 1976]

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Why reproducibility?

- Falsifiability and community consensus (Scientific Method)
- Building upon the work of others
 - Typical first step: reproduce previous results
- Quality control and error correction
 - *"Most published research findings are false"* ⇒ Reproducibility crisis
 - Yet, very few corrections (Errata) published
 - \Rightarrow Lack of reproducibility studies?

Reproducibility crisis in EC?

- "Most published research findings are false" [Ioannidis, 2005]
 Nature survey of 1 500 researchers: [Baker, 2016] +70% failed to reproduce another researcher's experiments +50% have failed to reproduce their own previous results
 Signs that the situation is CS is no better [Cockburn et al., 2020]
 In EC:
 - - Very few published cases
 - But no reason to think EC is special

[Sörensen et al., 2017]

What is Reproducibility?

• No consensus in terminology

[Claerbout & Karrenbach, 1992] [Plesser, 2018]

• ACM distinguishes between:



Repeatability, Reproducibility and Replicability

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

 López-Ibáñez, Branke, and Paquete [2021] define the terms more precisely and distinguish between: *Repeatability, Reproducibility, Replicability* and *Generalisability*

Terminology

Artifact

"A digital object that was either created by the authors to be used as part of the study or generated by the experiment itself"

[ACM, 2020]

algorithm implementations, benchmark instances, data pre/post-processing scripts, ...

Measurement [López-Ibáñez, Branke, and Paquete, 2021] "data that results from an experiment"

- measures of quality, computational effort, etc.
- NOT summary statistics

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization





ACM's Terminology



Association for Computing Machinery

Repeatability (Same team, same experimental setup)

Reproducibility (Different team, same experimental setup)

Replicability (Different team, different experimental setup)

•• The measurement can be obtained with stated precision by a different team, a different measuring system, in a different location on multiple trials.[...]

[A]n independent group can obtain the same result using artifacts which they develop completely independently.

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Dimensions of reproducibility

- *Artifacts:* Re-use of the original artifacts should allow to repeat the exact same experiments as described in the original publication
- Random factor:
 - the experiment only evaluates a random sample
 - the experimental claim applies to range or distribution
 - Random seeds

• Fixed factor:

- the experiment only evaluates specific chosen values
- experimental claim only supported for those specific values
- parameter settings, benchmark problems, computational budget ... unless randomized

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Proposed terminology [López-Ibáñez, Branke, and Paquete, 2021] of reproducibility studies

Label	Artifacts	Random factors	Fixed factors	Comment
Repeatability	Original	Original	Original	Exactly repeat the original experiment, generating precisely the same results.
Reproducibility	Original	New	Original	Test whether the original results were dependent on specific values of random factors and, hence, only a statistical anomaly.
Replicability	New	New	Original	Test whether it is possible to independently reach the same conclusion without relying on original artifacts.
Generalisability	Original or New	New	New	Test whether the conclusion extends beyond the experimental setup of the original paper. When new artifacts are used, generalisability should come after a replicability study.
Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimiz				

Obstacles to reproducibility

- Cultural obstacles
- Technical obstacles

Cultural Obstacles

- Disincentives to publish artifacts
 - **X** Additional effort \Rightarrow Fewer papers
 - **X** Error detection \Rightarrow Rejection / retraction
 - X Not *required* by journals *pre*-publication
- Difficulty to publish reproducibility studies
 - X More effort than a new algorithm / survey
 - X Low chances of publication
 - $\pmb{\mathsf{X}}$ Biases against both negative results and corrections
- Insufficient description
 - **X** No or bad artifacts \Rightarrow reproducibility impossible
 - **X** Paper often not enough for replication
 - "Obsolete" and correct code is **better** than no / incorrect code

Part II

Guidelines and Tools

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Technical Obstacles

- X Restrictive licenses, privacy and commercially sensitive data, etc
- **X** Binary only artifacts (no source code):
 - "Obsolete" <u>source code</u> is <u>better</u> than "working" <u>black-box executable</u>
- Unreproducible or unreplicable computational environment any difference may distort CPU-time, RNG, floating-point, etc
- Prohibitive or unavailable computational resources
 Years of CPU time or specific hardware (GPUs)
- Verification of artifacts
 - \mathbf{X} Manual verification \Rightarrow Tremendous effort
 - **X** Lack of re-implementations \Rightarrow Error propagation [Brockhoff, 2015]

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Ensuring reproducibility: Artifacts

- Permanently accessible (ACM badge artifact available):
 - X Personal / research group webpages repositories
 - Git tag / SHA commit: https://github.com/NEO-Research-Group/irace-sumo/ tree/62304739940199b3326cf8b34837c540cad6a68d
- Figshare (https://www.figshare.com), Zenodo (https://www.zenodo.org)
 Open Science Foundation (https://www.osf.io)
- Complete:
 - ✓ All source code and input data: Pre-processing code, Algorithm code, Analysis code and Presentation code
- Step-by-step documentation, flexible reproduction scripts and raw intermediate data
- Decision vectors (actual solutions), testsuite, independent solution checkers, ...
- Useful: Open-source license and open-data formats

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Zenodo Example	https://doi.org/10.5281/zenodo.4500973				
Zenodo Starch	Q Upload Communities		+9 Log in Cor Sign up		
February 4, 2021	Conference paper Open Access				
Unbalanced Mallows Mode Expensive Black-Box Perm	26 views See more of	3 Letails			
🕒 Irurozki, Ekhiñe; 🕜 López-Ibáñez, Manuel					
Reproducible Artifacts for the paper:					
Ekhine Irurozki and Manuel López-Ibáñez. Unbalanced Mallows Mode Permutation Problems. In Genetic and Evolutionary Computation Co ACM, New York, NY, USA, 9 pages. https://doi.org/10.1145/3449639.3	Indexed in				
Expensive black-box combinatorial optimization problems arise in pra means of a simulator or a real-world experiment. Since each fitness ar only a limited number of evaluations is possible, typically several orde problems. In this scenario, classical optimization methods such as m eachd. It is the aeritement of Devaluation and the second secon	Open	AIRE			
udertin m the commonscience, bayesian organization, in particular us under these conditions. Much les research is available in the combin UMM, an estimation-of-distribution (EDA) algorithm based on a Malion aggregation (uBorda). Experimental results on black-box versions of L sometimes surpass, the solutions obtained by CEGO, a Bayesian optim Moreover, the computational complexity of UMM increases linearly wi	Publication date: February 4, 2021 DOI: DOI: 10.5281/zenodo.4500	974			
permutation size.		Keyword(s): Combinatorial optimization Expensive black-box optimizat	Bayesian optimization		
Preview	*	Estimation of distribution algor	ithms		
🗈 umm zip	×	Creative Commons Attr	ibution 4.0 International		
The previewer is not showing all the files					
🗅 0-set up.sh					

Jupyter Python notebook https://doi.org/10.5281/zenodo.4500973



Ensuring reproducibility: Detailed experimental conditions

- X Many results are sensitive to computational platform:
 - CPU speed, cache sizes, floating-point arithmetic, library bugs (Linux's pow() bug 13932), ...
 - Virtual machines.
 - containers docker platforms 🔇 CODE OCEAN
 - **COSE**

Luís Paquete, Manuel López-Ibáñez

✓ Provide *calibration benchmark* running times

X Hidden / unfair parameter tuning:

- Report: Effort, Domains, Training problem instances
- Parameter tuning procedure should be reproducible:
 - ✓ Design of Experiments
 - [Montgomery, 2012] ✓ Automatic configuration tools, e.g., irace [López-Ibáñez et al., 2016]

Replicability and Reproducibility in Evolutionary Optimizatio

A checklist for reproducibility (1): Artifacts

- □ Step-by-step documentation to reproduce the experiment AND analysis
- □ All source code
- □ All input data: problem instances, random seeds, ...
- □ Analysis and presentation scripts
- □ Raw generated data (objective and decision vectors)
- □ Testsuite and solution checker
- □ Computation time calibration code and running times
- □ Open-data formats (✓ CSV ✓ MySQL ¥ Excel, ¥ Oracle, etc.)
- □ Open-source license (reading, distributing, running and reusing)
- □ Permanent link / DOI to specific version
- □ Long-term (permanently) accessible repository Personal website

GitHub repo

A checklist for reproducibility (2) Report / Document

- Relevant hardware details (CPU details, memory / cache sizes)
 Provide a container (e.g., Docker)
 - □ Provide link to virtual platform (e.g., Code Ocean)
 - □ Provide reviewer access to special hardware (e.g., GPUs)
- Precise versions of any additional software, packages, simulators, compilers / interpreters, and OS
- □ (Hyper-)parameters, including types and domains
- □ Parameter tuning process (also reproducible)
- □ Separate problem instances for development/tuning and for benchmarking / hypothesis testing
- □ Confidence intervals (or p-values), size effects

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Encouraging reproducibility efforts

- Journals should *require* that artifacts are provided to reviewers
 - Mathematical Programming Computation (Springer) requires source code
- (Ideally) Journals adopt the *Transparency and Openness Promotion (TOP)* guidelines: [Nosek et al., 2015; Stodden et al., 2016]
 - Reproducibility checks
 - Independent replication
 - ... among other requirements before publication
- ACM badges provide a way to recognise different degrees of reproducibility

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Reproducibility at ACM TELO

- Authors request a reproducibility level for their work
- A member of the reproducibility board at ACM TELO checks whether the work is at the requested level and may interact with the authors to improve it in terms of reproducibility
- Once approved, a *reproducibility badge* is provided



ACM Badges



Artifacts are permanently available for retrieval

ACM Badges



Artifacts evaluated Functional



Artifacts evaluated Reusable Same as above plus careful documentation in order to allow reusability

ACM Badges



Results validated Reproduced



Results validated Replicated The main results are obtained using the author-supplied artifacts

The main results are obtained without the use of author-supplied artifacts (not yet at ACM TELO)

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

The artifacts are well

requirements are met

documented and (functional)

Conclusions

- Reproducibility implies a *cultural shift*
- Authors, journals editors, conference organisers and funding institutes should be active players on promoting this shift
- Extra effort is rewarded with *faster scientific progress* and *higher reputation* for the field as a whole
- But is it such amount of effort? It only implies adopting good practices of verification, documentation and reporting
 - At GECCO, in the near future: "where is your code?"

Acknowledgments

The tutorial has benefited from collaborations and discussions with our colleagues:

Juergen Branke, Carola Doerr, Mike Preuss



The research leading to the results presented here has received funding from diverse projects:

- National funds through the FCT Foundation for Science and Technology, I.P. within the scope of the project CISUC – UID/CEC/00326/2020.
- M. López-Ibáñez is a "Beatriz Galindo" Senior Distinguished Researcher (BEAGAL 18/00053) funded by the Spanish Ministry of Science and Innovation (MICINN).

References I

- ACM. Artifact review and badging version 1.1. https://www.acm.org/publications/policies/artifact-review-and-badging-current, Aug. 2020.
- M. Baker. Is there a reproducibility crisis? Nature, 533:452-454, 2016
- D. Brockhoff. A bug in the multiobjective optimizer IBEA: Salutary lessons for code release and a performance re-assessment. In A. Gaspar-Cunha, C. H. Antunes, and C. A. Coello Coello, editors, *Evolutionary Multi-criterion Optimization, EMO 2015 Part I,* volume 9018 of *Lecture Notes in Computer Science*, pages 187–201. Springer, Heidelberg, Germany, 2015. doi: 10.1007/978-3-319-15934-8.13.
- J. Claerbout and M. Karrenbach. Electronic documents give reproducible research a new meaning. In SEG Technical Program Expanded Abstracts 1992, pages 601–604. Society of Exploration Geophysicists, 1992. doi: 10.1190/1.1822162.
- A. Cockburn, P. Dragicevic, L. Besançon, and C. Gutwin. Threats of a replication crisis in empirical computer science. Communications of the ACM, 63(8):70–79, July 2020. doi: 10.1145/3360311.
- J. P. A. loannidis. Why most published research findings are false. PLoS Medicine, 2(8):e124, 2005. doi: 10.1371/journal.pmed.0020124.
- M. López-Ibáñez, J. Dubois-Lacoste, L. Pérez Cáceres, T. Stützle, and M. Birattari. The irace package: Iterated racing for automatic algorithm configuration. Operations Research Perspectives, 3:43–58, 2016. doi: 10.1016/j.orp.2016.09.002.
- M. López-Ibáñez, J. Branke, and L. Paquete. Reproducibility in evolutionary computation. Arxiv preprint arXiv:20102.03380 [cs.Al], 2021. URL https://arxiv.org/abs/2102.03380.
- D. C. Montgomery. Design and Analysis of Experiments. John Wiley & Sons, New York, NY, 8th edition, 2012.
- Y. Nagata and S. Kobayashi. Edge assembly crossover: A high-power genetic algorithm for the traveling salesman problem. In T. Bäck, editor, ICGA, pages 450–457. Morgan Kaufmann Publishers, San Francisco, CA, 1997.
- A. Newell and H. A. Simon. Computer science as empirical inquiry: Symbols and search. Communications of the ACM, 19(3): 113–126, Mar. 1976. ISSN 0001-0782. doi: 10.1145/360018.360022.
- B. A. Nosek, G. Alter, G. C. Banks, D. Borsboom, S. D. Bowman, S. J. Breckler, S. Buck, C. D. Chambers, G. Chin, G. Christensen, M. Contestabile, A. Dafoe, E. Eich, J. Freese, R. Glennerster, D. Gordi, D. P. Green, B. Hesse, M. Humphreys, J. Ishiyama, D. Karian, A. Kraut, A. Lupia, P. Mabry, T. Madon, N. Malhotra, E. Mayo-Wilson, M. McNutt, E. Miguel, E. L. Paluck, U. Simonshn, C. Soderberg, B. A. Spellman, J. Turitto, G. VandenBos, S. Vazire, E. J. Wagenmakers, R. Wilson, and T. Yarkoni. Promoting an open research culture. Science, 348(6242):1422-1425, June 2015. doi: 10.1126/Science.ab2274.

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

References II

- H. E. Plesser. Reproducibility vs. replicability: A brief history of a confused terminology. Frontiers in Neuroinformatics, 11, Jan. 2018. doi: 10.3389/fninf.2017.00076.
- K. Sörensen, F. Arnold, and D. Palhazi Cuervo. A critical analysis of the "improved clarke and wright savings algorithm". International Transactions in Operational Research, 26(1):54–63, 2017. doi: 10.1111/itor.12443.
- V. Stodden, M. McNutt, D. H. Bailey, E. Deelman, Y. Gil, B. Hanson, M. A. Heroux, J. P. A. Ioannidis, and M. Taufer. Enhancing reproducibility for computational methods. *Science*, 354(6317):1240–1241, Dec. 2016. doi: 10.1126/science.aah6168.
- P. Wegner. Research paradigms in computer science. In ICSE'76: Proceedings of the 2nd international conference on Software engineering, pages 322–330, Oct. 1976.

Luís Paquete, Manuel López-Ibáñez Replicability and Reproducibility in Evolutionary Optimization

Replicability and Reproducibility in Evolutionary Optimization

Luís Paquete

paquete@dei.uc.pt http://www.uc.pt/go/paquete

University of Coimbra, Portugal



Manuel López-Ibáñez

manuel.lopez-ibanez@uma.es http://lopez-ibanez.eu

University of Málaga, Spain



