

TarGene: A Nextflow pipeline for the estimation of genetic effects on human traits via semi-parametric methods.

Olivier Labayle^{1,2}, Breeshey Roskams-Hieter^{1,2}, Joshua Slaughter^{1,2}, Kelsey Tetley-Campbell^{1,2}, Mark J. van der Laan⁴, Chris P. Ponting³, Ava Khamseh^{1,2,4}, and Sjoerd Viktor Beentjes^{1,3,4}

¹ MRC Human Genetics Unit, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh EH4 2XU, United Kingdom. ² School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, United Kingdom ³ School of Mathematics and Maxwell Institute, University of Edinburgh, Edinburgh EH9 3FD, United Kingdom ⁴ Division of Biostatistics, University of California, Berkeley, CA, USA

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Genetic variations are the foundation of biological diversity, they play a crucial role in the adaptability, survival, and evolution of populations. Discovering which and how genetic variations affect human traits is an ongoing challenge with applications in healthcare and medicine. In some cases, genetic variations have an obvious effect because they change the coding sequence of a gene and thus its function. In the vast majority of cases however, variations occur in places of unknown function and could impact human traits or disease mechanisms in complex ways. TarGene is a Nextflow pipeline leveraging highly flexible machine-learning methods and semi-parametric estimation theory to capture these complex genetic dependencies including higher-order interactions.

Statement of Need

All currently existing software for the estimation of genetic effects are based on parametric distributions, additionally assuming linearity of the relationship between variants and traits (Purcell et al., 2007, pp. [yang2011gcta](#), [loh2018mixed](#), [zhou2018efficiently](#)). If these assumptions are violated, the reported effect sizes will be biased and error rates inflated. In particular, this can lead to inflated false discovery rates and suboptimal resources allocation. Some recently published software also account for more complex relationships but do not offer the full modelling flexibility provided by TarGene. REGENIE has the benefit to fit a whole-genome model for each phenotype of interest but still assumes linearity and normality (Mbatchou et al., 2021). DeepNull is a semi-parametric method which models non-linear covariate effects but also assumes genetic effects to be linear and does not allow complex interactions between covariates and genetic variants (McCaw et al., 2022). KnockoffGWAS (Sesia et al., 2021), is non-parametric but does not estimate effect sizes, instead it aims at controlling the false discovery rate in genome-wide association studies. In comparison, TarGene is the only method able to model arbitrarily complex genetic effects while preserving the validity of statistical inferences. It does so by leveraging Targeted Learning (Van der Laan et al., 2011), a framework combining methods from causal inference, machine-learning and semi-parametric statistical theory. Succinctly, the estimation process works as follows. In an first step, flexible machine-learning algorithms are fitted to the data, hence minimizing an appropriate loss function (e.g., negative log-likelihood). A second step, known as the targeting step, reduces the estimation bias in a theoretically optimal way.

Features

TarGene is a fully featured command-line software, which can be run as follow:

```
nextflow run https://github.com/TARGENE/targene-pipeline/ \
  -r TARGENE_VERSION \
  -c CONFIG_FILE \
  -resume
```

where the CONFIG_FILE provides the list of problem specific parameters (data, arguments, options). Below we list some important features of TarGene, the following CONFIG_FILE will serve as a running example.

```
params {
  ESTIMANDS_CONFIG = "gwas_config.yaml"
  ESTIMATORS_CONFIG = "wtmle--tunedxgboost"

  // UK-Biobank specific parameters
  BED_FILES = "unphased_bed/ukb_chr{1,2,3}.{bed,bim,fam}"
  UKB_CONFIG = "ukbconfig_gwas.yaml"
  TRAITS_DATASET = "dataset.csv"
}
```

For detailed explanations, please refer to the online [documentation](#).

Scalability

Machine-learning methods are computationally intensive, however statistical genetics analyses need to scale to hundreds of thousands of variants and thousands of traits. For this reason, TarGene leverages Nextflow ([Di Tommaso et al., 2017](#)), a pipeline management system that can parallelize independent estimation tasks across HPC platforms.

Databases

TarGene works with standard formats, plink .bed and .bgen formats for genotypes, .csv or .arrow format for human traits. Furthermore, TarGene has direct support for two large scale biomedical databases, the UK-Biobank ([Bycroft et al., 2018](#)) and the All of Us cohort ([Us Research Program Investigators, 2019](#)). The example considers the UK-Biobank for which genotypes and traits are provided via BED_FILES and TRAITS_DATASET respectively. Because the UK-Biobank has a non-standard format, the UKB_CONFIG provides traits definition rules. The following is an illustration for BMI, but the default is to consider all 766 traits as defined by the geneAtlas ([Canela-Xandri et al., 2018](#)).

```
traits:
  - fields:
    - "21001"
  phenotypes:
    - name: "Body mass index (BMI)"
```

Study Designs

TarGene supports traditional study designs in population genetics, that is, genome-wide association studies (GWAS) and phenome-wide association studies (PheWAS). Because TarGene has a focus on complex effects, higher-order interactions (e.g. gene-gene-... or gene-environment-...) can also be investigated.

The study design is specified in the ESTIMANDS_CONFIG YAML file. For a routine GWAS the content of this file can be as simple as:

87 type: gwas

88 Estimators

89 Semi-parametric estimators exist in multiple flavors, all with different properties. In TarGene
90 we default to using Targeted Maximum-Likelihood Estimation (Van der Laan & Rose, 2018)
91 and XGboost (Chen & Guestrin, 2016) as the machine-learning model. This is because this
92 was the best performing estimator in simulations for a variety of tasks. But if computational
93 restrictions exist, tradeoffs can be made and simpler models can be used.

94 Acknowledgements

95 This work was supported by the United Kingdom Research and Innovation (grant
96 EP/S02431X/1), UKRI Centre for Doctoral Training in Biomedical AI at the University of
97 Edinburgh, School of Informatics.

98 References

- 99 Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K., Motyer, A., Vukcevic,
100 D., Delaneau, O., O'Connell, J., & others. (2018). The UK biobank resource with deep
101 phenotyping and genomic data. *Nature*, 562(7726), 203–209.
- 102 Canela-Xandri, O., Rawlik, K., & Tenesa, A. (2018). An atlas of genetic associations in UK
103 biobank. *Nature Genetics*, 50(11), 1593–1599.
- 104 Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. *Proceedings of
105 the 22nd Acm Sigkdd International Conference on Knowledge Discovery and Data Mining*,
106 785–794.
- 107 Di Tommaso, P., Chatzou, M., Floden, E. W., Barja, P. P., Palumbo, E., & Notredame, C.
108 (2017). Nextflow enables reproducible computational workflows. *Nature Biotechnology*,
109 35(4), 316–319.
- 110 Mbatchou, J., Barnard, L., Backman, J., Marcketta, A., Kosmicki, J. A., Ziyatdinov, A.,
111 Benner, C., O'Dushlaine, C., Barber, M., Boutkov, B., & others. (2021). Computationally
112 efficient whole-genome regression for quantitative and binary traits. *Nature Genetics*, 53(7),
113 1097–1103.
- 114 McCaw, Z. R., Colthurst, T., Yun, T., Furlotte, N. A., Carroll, A., Alipanahi, B., McLean, C.
115 Y., & Hormozdiari, F. (2022). DeepNull models non-linear covariate effects to improve
116 phenotypic prediction and association power. *Nature Communications*, 13(1), 241.
- 117 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., Maller,
118 J., Sklar, P., De Bakker, P. I., Daly, M. J., & others. (2007). PLINK: A tool set for
119 whole-genome association and population-based linkage analyses. *The American Journal
120 of Human Genetics*, 81(3), 559–575.
- 121 Sesia, M., Bates, S., Candès, E., Marchini, J., & Sabatti, C. (2021). False discovery rate
122 control in genome-wide association studies with population structure. *Proceedings of the
123 National Academy of Sciences*, 118(40), e2105841118.
- 124 Us Research Program Investigators, A. of. (2019). The “all of us” research program. *New
125 England Journal of Medicine*, 381(7), 668–676.
- 126 Van der Laan, M. J., & Rose, S. (2018). *Targeted learning in data science*. Springer.
- 127 Van der Laan, M. J., Rose, S., & others. (2011). *Targeted learning: Causal inference for
128 observational and experimental data* (Vol. 4). Springer.