

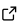
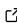
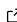
1 PhenoFeatureFinder: a python package for linking 2 developmental phenotypes to omics features

3 Lissy-Anne M. Denkers ¹, Marc D. Galland ², Annabel Dekker³, Valerio
4 Bianchi ³, and Petra M. Bleeker ¹

5 ¹ University of Amsterdam, Department of Plant Physiology, Green Life Science Research Theme,
6 Swammerdam Institute for Life Sciences, Amsterdam, The Netherlands ² INRAE, Institute of Genetics,
7 Environment and Plant Protection (IGEPP—Joint Research Unit 1349), Le Rheu, France ³ Enza Zaden
8 R&D B.V., BTR-BM Bioinformatics, Enkhuizen, The Netherlands

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

Software

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

Editor: Julia Romanowska 

Submitted: 05 June 2024

Published: unpublished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](https://creativecommons.org/licenses/by/4.0/)).

9 Summary

10 Plants interact with their (a)biotic environment through a range of specialised metabolites.
11 They deal with pathogens and pest attack through constitutive or inducible production of
12 specialised metabolites or other defence molecules ([Erb & Kliebenstein, 2020](#); [García-Olmedo
13 et al., 1998](#)). High-throughput “-omics” tools including (untargeted) metabolomics have been
14 successfully implemented in plant biology ([Dalio et al., 2021](#)), but the accompanying resistance
15 phenotyping often lacks in robustness ([Song et al., 2021](#)). Resistance is often not a binary
16 trait, but quantitative in nature. To identify features such as defence metabolites involved in a
17 resistance phenotype, we developed a pipeline that includes phenotypic analysis, preprocessing
18 and visualisation of the metabolomics data, and feature prediction through a Machine Learning
19 approach.

20 Proliferation of an insect population is affected by various factors, including the chemical
21 composition of the host, and/or the environment ([Ma et al., 2022](#)). In particular, host resistance
22 via hampered larval development is noteworthy, because reducing the speed at which larvae
23 reach the adult stage and produce offspring negatively affects pest-population development
24 ([Maharijaya et al., 2019](#); [Muema et al., 2016](#); [Vengateswari et al., 2022](#)). However, evaluating
25 larval development results in a complex dataset that is challenging to process. Developmental
26 success is based on the number of larvae throughout various larval stages, as well as on the
27 speed of development.

28 To identify underlying mechanisms of resistance, the chemical or molecular composition of a
29 plant can be investigated. Proteins and metabolites are commonly analysed through untargeted
30 Mass-Spectrometry, yielding exhaustive profiles generally consisting of many thousands of
31 unannotated features. Often such data displays sparsity, i.e. missing values between datasets,
32 and a low sample-to-feature ratio, adding to the complexity of the analysis ([Kortbeek et al.,
33 2021](#); [Liebal et al., 2020](#)). Tree-based Machine-Learning algorithms (e.g., random forest) are
34 suitable for the analysis of, and feature selection from, untargeted data ([Liebal et al., 2020](#))
35 computing the contribution of each feature in the phenotypic classification.

36 PhenoFeatureFinder is designed to facilitate the different analyses mentioned above in one
37 pipeline. It can be used for 1) evaluation and visualisation of pest performance over multiple
38 stages and between groups (treatments, genotypes), 2) pre-processing of the omics data, and
39 3) prediction of features that explain the phenotypic classification. To facilitate usability, each
40 step in the pipeline can also be performed independently, hence has been assigned a class
41 in the package ([Figure 1](#)). Also, although we focus on insect development and the selection
42 of metabolic features causal to the observed phenotype, different input data with a similar
43 structure could be used. PhenoFeatureFinder was developed initially for metabolomics data,

44 users can evaluate its fit applying other types of omics data.

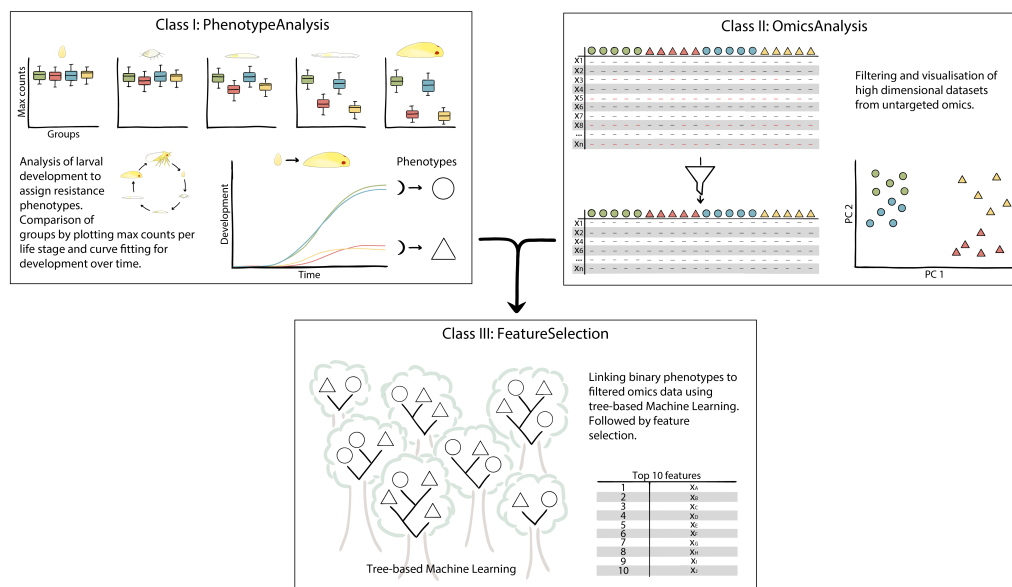


Figure 1: Overview of the package, consisting of three classes that can be used separately or as a workflow. Class 1: analysing and visualising the phenotype, Class 2: preprocessing and visualising omics datasets, and Class 3: feature selection through a Machine Learning approach.

45 Statement of need

46 Class I: PhenotypeAnalysis

47 A binary classification of plants into “resistant” or “susceptible” helps to extract relevant
48 features especially when threshold effects or sparsity (presence/absence) effects are at play.
49 Here we firstly assess performance over different developmental stages of larvae on different
50 host plants. The number of individuals in each stage at a given time is recorded. When plotted,
51 the cumulative data of these bioassays resemble a growth- or dose-response curve that can
52 be used to manually assign a binary phenotype (e.g., resistant/non-resistant), a resistance
53 classification used as input for FeatureSelection (Class 3).

54 A package named `drc` is available in R for fitting dose-response curves (Ritz et al., 2015),
55 offering an extensive and versatile set of functionalities. However, for the purposes described
56 here `drc` poses some limitations, such as the options for custom pre-processing and analyses
57 of multiple experimental groups simultaneously. Here we implemented pre-processing steps
58 and aimed to decrease the amount of coding needed to obtain a fitted development curve.
59 To account for missing data when individuals that reached the final developmental stage are
60 removed from the experiment, we implemented an automated correction step. The count
61 data can be transformed to cumulative data to analyse the maximum of individuals that reach
62 each of the developmental stages. Next, the time to reach a specific stage can be compared
63 between treatments by fitting a 3-parameter log-logistic curve (Muse et al., 2021; Seefeldt et
64 al., 1995; Vliet & Ritz, 2013) to the cumulative data for each treatment, with the function:

$$f(x) = \frac{m}{1 + \exp(s \times (\log(x) - \log(e_{50})))}$$

65 where x is time, m is the upper limit (or maximum of individuals that developed to the stage of
66 interest), s is the slope of the linear part of the curve and e_{50} is the EmT50 (the timepoint at

67 which 50% of the individuals have developed to the stage of interest). We added the possibility
68 to compare performance between treatments by fitting a curve with the function:

$$f(x) = \frac{a \times \frac{s}{m} \times \left(\frac{x}{m}\right)^{s-1}}{1 + \left(\frac{x}{m}\right)^s}$$

69 Here, x is time, a the area under the curve, s is the shape of the curve and m the median
70 time point. Both functions output a table with the model parameters, confidence intervals
71 and the model fit, together with a plot displaying the observed data and the fitted model. For
72 both functions it is possible to predict the potential maximum beyond the final experimental
73 measurements.

74 Class II: OmicsAnalysis

75 Untargeted omics results in large datasets that tend to contain background noise and unreliable
76 features. To clean the data, multiple filtering methods are implemented in the `OmicsAnalysis`
77 class, including the removal of contaminants present in blank samples, filtering to decrease
78 sparsity and other quality control steps. The structure of the data can subsequently be
79 visualised with a PCA and an UpSet plot.

80 Class III: FeatureSelection

81 Combining the output of Classes 1 and 2, i.e. the binary phenotype classification and the
82 tidied untargeted metabolomics, `FeatureSelection` is set up to predict features that can
83 explain the phenotypic observation under study. This part of the pipeline was built as a
84 wrapper around the Python libraries `scikit-learn` and `TPOT` (Olson et al., 2016; Pedregosa
85 et al., 2011). The `FeatureSelection` wrapper is designed to select optimal pipelines for
86 data preprocessing and identification of the most suitable Machine Learning model. One
87 characteristic of metabolomics data is strongly correlated features (linear dependencies between
88 variables) that make it difficult to extract individual feature importance. Therefore, this method
89 implements a PCA as dimensionality reduction method before searching for the best fitting
90 pipeline. Finally, the importance of the Principal Components and their most related features
91 (high loadings) can be retrieved to select features with predicted importance to the phenotypic
92 classification.

93 Acknowledgements

94 The Amsterdam Data Science Centre is acknowledged for their input and Frans van der Kloet
95 for his statistical support. Gea van der Lee and Piet Verdonshot are kindly thanked for
96 providing the data from (Lee et al., 2020) used in the example in GitHub.

97 References

- 98 Dalio, R. J. D., Litholdo, C. G., Arena, G., Magalhães, D., & Machado, M. A. (2021).
99 *Contribution of omics and systems biology to plant biotechnology* (pp. 171–188). https://doi.org/10.1007/978-3-030-80352-0_10
- 100
- 101 Erb, M., & Kliebenstein, D. J. (2020). Plant secondary metabolites as defenses, regulators,
102 and primary metabolites: The blurred functional trichotomy. *Plant Physiology*, 184, 39–52.
103 <https://doi.org/10.1104/pp.20.00433>
- 104 García-Olmedo, F., Molina, A., Alamillo, J. M., & Rodríguez-Palenzuela, P. (1998).
105 Plant defense peptides. *Biopolymers*, 47, 479–491. [https://doi.org/10.1002/\(SICI\)1097-0282\(1998\)47:6%3C479::AID-BIP6%3E3.0.CO;2-K](https://doi.org/10.1002/(SICI)1097-0282(1998)47:6%3C479::AID-BIP6%3E3.0.CO;2-K)
- 106

- 107 Kortbeek, R. W. J., Galland, M. D., Muras, A., Kloet, F. M. van der, André, B., Heilijgers,
108 M., Hijum, S. A. F. T. van, Haring, M. A., Schuurink, R. C., & Bleeker, P. M. (2021).
109 Natural variation in wild tomato trichomes; selecting metabolites that contribute to
110 insect resistance using a random forest approach. *BMC Plant Biology*, *21*, 315. <https://doi.org/10.1186/s12870-021-03070-x>
111
- 112 Lee, G. H. van der, Kraak, M. H. S., Verdonschot, R. C. M., & Verdonschot, P. F. M.
113 (2020). Persist or perish: Critical life stages determine the sensitivity of invertebrates to
114 disturbances. *Aquatic Sciences*, *82*. <https://doi.org/10.1007/s00027-020-0698-0>
- 115 Liebal, U. W., Phan, A. N. T., Sudhakar, M., Raman, K., & Blank, L. M. (2020). Machine
116 learning applications for mass spectrometry-based metabolomics. *Metabolites*, *10*, 243.
117 <https://doi.org/10.3390/metabo10060243>
- 118 Ma, K., Tang, Q., Liang, P., Li, J., & Gao, X. (2022). A sublethal concentration of afidopyropen
119 suppresses the population growth of the cotton aphid, *aphis gossypii* glover (hemiptera:
120 aphididae). *Journal of Integrative Agriculture*, *21*, 2055–2064. [https://doi.org/10.1016/S2095-3119\(21\)63714-0](https://doi.org/10.1016/S2095-3119(21)63714-0)
121
- 122 Maharijaya, A., Vosman, B., Pelgrom, K., Wahyuni, Y., Vos, R. C. H. de, & Voorrips, R. E.
123 (2019). Genetic variation in phytochemicals in leaves of pepper (*capsicum*) in relation
124 to thrips resistance. *Arthropod-Plant Interactions*, *13*, 1–9. <https://doi.org/10.1007/s11829-018-9628-7>
125
- 126 Muema, J. M., Bargul, J. L., Nyanjom, S. G., Mutunga, J. M., & Njeru, S. N. (2016).
127 Potential of *camellia sinensis* proanthocyanidins-rich fraction for controlling malaria mosquito
128 populations through disruption of larval development. *Parasites & Vectors*, *9*, 512. <https://doi.org/10.1186/s13071-016-1789-6>
129
- 130 Muse, A. H., Mwalili, S. M., & Ngesa, O. (2021). On the log-logistic distribution and its
131 generalizations: A survey. *International Journal of Statistics and Probability*, *10*, 93.
132 <https://doi.org/10.5539/ijsp.v10n3p93>
- 133 Olson, R. S., Bartley, N., Urbanowicz, R. J., & Moore, J. H. (2016). Evaluation of a tree-based
134 pipeline optimization tool for automating data science. *Proceedings of the Genetic and
135 Evolutionary Computation Conference 2016*, 485–492. <https://doi.org/10.1145/2908812.2908918>
136
- 137 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
138 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D.,
139 Brucher, M., Perrot, M., & Duchesnay, É. (2011). Scikit-learn: Machine learning in python.
140 *Journal of Machine Learning Research*, *12*, 2825–2830. <https://doi.org/10.5555/1953048.2078195>
141
- 142 Ritz, C., Baty, F., Streibig, J. C., & Gerhard, D. (2015). Dose-response analysis using r. *PLoS
143 ONE*, *10*, 1–13. <https://doi.org/10.1371/journal.pone.0146021>
- 144 Seefeldt, S. S., Jensen, J. E., & Fuerst, E. P. (1995). Log-logistic analysis of herbicide
145 dose-response relationships. *Weed Technology*, *9*, 218–227. <https://www.jstor.org/stable/3987736>
146
- 147 Song, P., Wang, J., Guo, X., Yang, W., & Zhao, C. (2021). High-throughput phenotyping:
148 Breaking through the bottleneck in future crop breeding. *The Crop Journal*, *9*, 633–645.
149 <https://doi.org/10.1016/j.cj.2021.03.015>
- 150 Vengateswari, G., Arunthirumeni, M., Shivaswamy, M. S., & Shivakumar, M. S. (2022). Effect
151 of host plants nutrients, antioxidants, and phytochemicals on growth, development, and
152 fecundity of *spodoptera litura* (fabricius) (lepidoptera: noctuidae). *International Journal of
153 Tropical Insect Science*, *42*, 3161–3173. <https://doi.org/10.1007/s42690-022-00868-6>
- 154 Vliet, L. van der, & Ritz, C. (2013). Statistics for analyzing ecotoxicity test data. In J.-F. Féraud

155 & C. Blaise (Eds.), *Encyclopedia of Aquatic Ecotoxicology* (pp. 1081–1095). Springer
156 Dordrecht. <https://doi.org/10.1007/978-94-007-5704-2>

DRAFT