



## GenForge — A Multi-Population GP Toolkit with SPFP

# GenForge User Manual & Tutorial

**Multi-Population Genetic Programming for Classification & Regression Tasks  
with Semantic Preserving Feature Partitioning (SPFP)**

***Mohammad Sadegh Khorshidi<sup>1</sup> , Navid Yazdanjue<sup>1</sup> ,***  
***Hassan Gharoun<sup>1</sup> , Mohammad Reza Nikoo<sup>2</sup> , Fang Chen<sup>1</sup> ,***  
***Amir H. Gandomi<sup>1,3</sup> ***

<sup>1</sup>*Faculty of Engineering & Information Technology, University of Technology Sydney, Ultimo 2007, Australia.*

<sup>2</sup>*Department of Civil and Architectural Engineering, Sultan Qaboos University, Muscat, Oman.*

<sup>3</sup>*University Research and Innovation Center (EKIK), Óbuda University, Budapest 1034, Hungary.*

*Correspondence:* msadegh.khorshidi.ak@gmail.com

October 28, 2025

Version 1.0



GenForge package emblem.

# Preface

This manual documents `GenForge`, a Python package for multi-population Genetic Programming that supports classification and regression tasks via `gpclassifier` and `gpregressor`, and semantic preserving feature partitioning via `SPFPPartitioner`. It covers installation, end-to-end examples, configuration and API references, plotting/reporting utilities, and design rationales.

**How to cite.** If you use `GenForge` or SPFP in academic work, please cite the GenForge repository and the following manuscripts:

- Khorshidi, M. S., Yazdanjue, N., Gharoun, H., Yazdani, D., Nikoo, M. R., Chen, F., & Gandomi, A. H. (2025). Semantic-Preserving Feature Partitioning for multi-view ensemble learning. *Information Fusion*, 122, 103152. DOI: [10.1016/j.inffus.2025.103152](https://doi.org/10.1016/j.inffus.2025.103152)
- Khorshidi, M. S., Yazdanjue, N., Gharoun, H., Nikoo, M. R., Chen, F., & Gandomi, A. H. (2025). Multi-population Ensemble Genetic Programming via Cooperative Coevolution and Multi-view Learning for Classification. *arXiv preprint arXiv:2509.19339*. DOI: [10.48550/arXiv.2509.19339](https://doi.org/10.48550/arXiv.2509.19339)
- Khorshidi, M. S., Yazdanjue, N., Gharoun, H., Nikoo, M. R., Chen, F., & Gandomi, A. H. (2025). From embeddings to equations: Genetic-programming surrogates for interpretable transformer classification. *arXiv preprint arXiv:2509.21341*. DOI: [10.48550/arXiv.2509.21341](https://doi.org/10.48550/arXiv.2509.21341)
- Khorshidi, M. S., Yazdanjue, N., Gharoun, H., Nikoo, M. R., Chen, F., & Gandomi, A. H. (2025). Domain-Informed Genetic Superposition Programming: A Case Study on SFRC Beams. *arXiv e-prints*, *arXiv-2509*. *arXiv preprint arXiv:2509.21341*. DOI: [10.48550/arXiv.2509.21355](https://doi.org/10.48550/arXiv.2509.21355)



# Contents

<b>Preface</b>	<b>iii</b>
<b>List of Figures</b>	<b>ix</b>
<b>List of Tables</b>	<b>xi</b>
<b>1 Introduction &amp; Quickstart</b>	<b>1</b>
1.1 What is GenForge? . . . . .	1
1.2 High-level architecture . . . . .	1
1.3 Installation . . . . .	1
1.4 First run . . . . .	2
1.4.1 Classification Task: <code>gpclassifier</code> . . . . .	2
1.4.2 Regression Task: <code>gpregressor</code> . . . . .	3
1.4.3 Feature partitioning: <code>SPFPPartitioner</code> . . . . .	4
1.5 Reproducibility notes . . . . .	5
<b>2 Installation and Environment Setup</b>	<b>7</b>
2.1 System Requirements . . . . .	7
2.2 Install from PyPI . . . . .	7
2.3 Install from Source . . . . .	8
2.4 Verifying the Installation . . . . .	8
2.5 Parallelism and Reproducibility . . . . .	9
2.6 Data Preparation . . . . .	9
2.7 Assets to Include for This Chapter . . . . .	9
<b>3 <code>gpclassifier</code> &amp; <code>gpregressor</code> Configurations</b>	<b>11</b>
3.1 Overview . . . . .	11
3.2 Type conventions used for <code>gpclassifier</code> . . . . .	12

3.3	Type conventions used for <code>gpregressor</code>	13
3.4	User Data	14
3.5	Run Control	14
3.6	Selection	16
3.7	Nodes: Ephemeral Constants and Function Set	17
3.8	Genetic Operators	17
3.9	Gene-Level Controls	18
3.10	Tree Constraints	18
3.11	Fitness and Objectives	18
3.12	Softmax Head (Specific to <code>ClassifierConfig</code> )	19
3.13	ElasticNetCV Head (Specific to <code>RegressorConfig</code> )	20
3.14	Resolved dictionaries	21
3.15	Common validation rules (summary)	21
3.16	Complete example	22
<b>4</b>	<b>SPFP: Mathematics and Algorithmic Details</b>	<b>25</b>
4.1	Notation and Discretisation	25
4.2	Per-Feature Scoring and Ranking	26
4.3	Group-Wise Forward Selection	26
4.4	Optional Backward Pruning	28
4.5	Forming Multiple Groups and Partition Policy	29
4.6	Parameter Summary	29
4.7	Complexity and Implementation Notes	30
4.8	Pseudocode	30
<b>5</b>	<b>gpclassifier: Legacy Configuration &amp; API</b>	<b>33</b>
5.1	Overview & Data Requirements for <code>gpclassifier</code>	33
5.2	Training Pipeline	33
5.3	Configuration Keys	34
5.3.1	User data & run metadata	34
5.3.2	Run control & parallelism	34
5.3.3	Function sets, tree/gene constraints	35
5.3.4	Selection & ensemble	35
5.3.5	Plotting & reporting	35

5.4	Using SPFP Partitions . . . . .	36
5.5	Training & Evaluation Workflow . . . . .	37
5.6	Outputs: Plots and HTML Reports . . . . .	38
5.7	Tips & Common Pitfalls . . . . .	38
<b>6</b>	<b>gpregressor: Legacy Configuration &amp; API</b>	<b>41</b>
6.1	Overview & Data Requirements for gpregressor . . . . .	41
6.2	Training Pipeline . . . . .	41
6.3	Configuration Keys . . . . .	42
6.3.1	User data & run metadata . . . . .	42
6.3.2	Run control & parallelism . . . . .	42
6.3.3	Function sets, tree/gene constraints . . . . .	43
6.3.4	Selection & ensemble . . . . .	43
6.3.5	Plotting & reporting . . . . .	43
6.4	Using SPFP Partitions . . . . .	44
6.5	Training & Evaluation Workflow . . . . .	45
6.6	Outputs: Plots and HTML Reports . . . . .	46
6.7	Tips & Common Pitfalls . . . . .	46
<b>7</b>	<b>SPFPPartitioner: Configuration &amp; API</b>	<b>49</b>
7.1	Overview . . . . .	49
7.2	Constructor Parameters . . . . .	49
7.3	Public Methods . . . . .	50
7.4	Minimal Usage . . . . .	50
7.5	Wiring Partitions to GP . . . . .	52
7.6	Notes & Gotchas . . . . .	53
<b>8</b>	<b>Plots &amp; HTML Reports</b>	<b>55</b>
8.1	Enabling and Saving Plots . . . . .	55
8.2	Fitness vs. Generation . . . . .	56
8.3	RankBest vs. Generation . . . . .	56
8.4	RankAll vs. Generation . . . . .	56
8.5	HTML Ensemble Report . . . . .	56
8.6	Interpreting the Ensemble . . . . .	57

8.7	Assets to Include for This Chapter	58
8.8	Troubleshooting	58
<b>9</b>	<b>Reproducibility &amp; Experiment Management</b>	<b>63</b>
9.1	Goals and Scope	63
9.2	Seeding and Determinism	63
9.3	Controlling BLAS & Threading	64
9.4	Data Splits and Consistency	65
9.5	Configuration Management	65
9.6	Run Naming and Artifact Layout	66
9.7	Logging Metrics and Comparing Runs	67
9.8	Strict vs. Practical Determinism	68
9.9	End-to-End Example	68
9.10	Troubleshooting	70
<b>10</b>	<b>Advanced Usage &amp; Extensibility</b>	<b>71</b>
10.1	Wiring SPFP Partitions to Multi-Population GP	71
10.2	Custom Primitive Function Sets	72
10.2.1	Selecting built-in tokens	72
10.2.2	Passing custom callables	73
10.3	Constants & Terminals	73
10.4	Tree & Gene Constraints	74
10.5	Selection & Ensemble Mixing	74
10.6	Parallel Fitness Evaluation	75
10.7	Saving, Logging, & Reproducibility	75
10.8	Performance Tuning Checklist	76
10.9	Minimal End-to-End Template	76
10.10	Troubleshooting Extensions	78

# List of Figures

1.1	GenForge data flow: (1) optional SPFP partitions features into $P$ groups ( <i>findex</i> ); (2) each group maps to a GP population; (3) a linear mixing head combines population outputs into the final prediction. . . . .	6
2.1	Example local project layout for GenForge experiments. . . . .	10
5.1	Example gpclassifier Fitness vs. Generation. . . . .	39
5.2	Rendered gpclassifier HTML report. . . . .	40
6.1	Example gpregressor Fitness vs. Generation. . . . .	47
6.2	Rendered gpregressor HTML report. . . . .	48
8.1	Fitness vs. Generation for a representative run. . . . .	59
8.2	RankBest vs. Generation (lower is better). . . . .	60
8.3	RankAll vs. Generation showing median and dispersion. . . . .	61
8.4	HTML report (example) summarising an ensemble row. . . . .	62
9.1	Example runs/ directory with plots, reports, and logs. . . . .	67



# List of Tables

2.1	Supported environments and core dependencies.	7
4.1	SPFP key parameters and recommended defaults (matching the reference implementation).	29
5.1	User data and metadata keys for <code>gpclassifier</code> .	34
5.2	Run control and parallel settings.	34
5.3	Structural configuration (function sets, depth, genes).	35
5.4	Selection and ensemble hyperparameters.	35
5.5	Plot and report toggles.	35
6.1	User data and metadata keys for <code>gpregressor</code> .	42
6.2	Run control and parallel settings.	42
6.3	Structural configuration (function sets, depth, genes).	43
6.4	Selection and ensemble hyperparameters.	43
6.5	Plot and report toggles.	43
7.1	Public constructor of <code>SPFPPartitioner</code> . All defaults mirror the reference implementation.	50
7.2	Core methods of <code>SPFPPartitioner</code> .	51
8.1	Common plot flags and output controls.	55
9.1	Useful environment variables for threading control.	64
9.2	Suggested minimal keys to track per run (extend as needed).	66



# 1 Introduction & Quickstart

## 1.1 What is GenForge?

GenForge is a multi-population Genetic Programming (GP) toolkit for classification and regression tasks. It couples (i) *functional diversity* via multiple co-evolving populations with (ii) *ensemble mixing* across populations, and augments the pipeline with *Semantic Preserving Feature Partitioning* (SPFP) to form feature groups that preserve information about the target. The package exposes three top-level modules:

- `gpclassifier` — multi-population GP for supervised classification task,
- `gpregressor` — multi-population GP for supervised regression task,
- `SPFPPartitioner` — a feature partitioner producing `num_groups` index sets (`findex`) to wire features to GP populations.

## 1.2 High-level architecture

At a high level, users (optionally) run `SPFPPartitioner` to obtain a partition  $\{\mathcal{I}_1, \dots, \mathcal{I}_P\}$  of input features, map each partition to a GP population via `gpclassifier` or `gpregressor`, then train a GP model whose ensemble head mixes per-population outputs. Figure 1.1 sketches the flow.

## 1.3 Installation

The package is available from its GitHub repository and can be installed from source.

**Listing 1.1:** Clone and install from GitHub.

```
1 git clone https://github.com/maisamkhorshidi/genforge.git  
2 cd genforge
```

```
3 pip install -e .
```

**Python compatibility.** GenForge supports Python 3.10+ (recommended). It is also known to run on Python 3.9.

**Dependencies.** Required: numpy, scipy, scikit-learn, joblib, pandas, matplotlib, seaborn. GenForge uses scikit-learn's ElasticNetCV for the fusion head,

**Versions.** Verified modern stack (Python  $\geq$ 3.10): joblib 1.5.2, numpy 2.3.4, pandas 2.3.3, scikit-learn 1.7.2, scipy 1.16.2, matplotlib 3.10.7, seaborn 0.13.2.

For Python 3.9, use a legacy stack, for example: numpy 1.26.x, scipy 1.11–1.13, scikit-learn 1.3–1.5, pandas 1.5–2.2, joblib  $\geq$ 1.3. This preserves API compatibility (e.g., ElasticNetCV) while meeting older interpreter constraints.

## 1.4 First run

### 1.4.1 Classification Task: gpclassifier

**Listing 1.2:** Minimal classification example with gpclassifier.

```
1 from genforge.gpclassifier import gpclassifier
2 import numpy as np
3
4 # Toy splits (NxD arrays and 1D labels)
5 X_train, y_train = np.load("X_train.npy"), np.load("y_train.npy")
6 X_val, y_val = np.load("X_val.npy"), np.load("y_val.npy")
7 X_test, y_test = np.load("X_test.npy"), np.load("y_test.npy")
8
9 # Optional: feature partitions (list of lists), one list per population
10 # findex = [[...], [...], ...]
11
12 params = {
13     "userdata_name": "demo_clf",
14     "userdata_xtrain": X_train, "userdata_ytrain": y_train,
```

```

15 "userdata_xval": X_val, "userdata_yval": y_val,
16 "userdata_xtest": X_test, "userdata_ytest": y_test,
17 # "userdata_pop_idx": findex, # (optional) wire features to populations
18 "runcontrol_pop_size": 64,
19 "runcontrol_generations": 50,
20 "runcontrol_useparallel": False,
21 "functions_name": [["times","minus","plus","divide"]],
22 "gene_max_genes": [5],
23 "tree_max_depth": [10],
24 "selection_elite_fraction": [0.05],
25 "selection_elite_fraction_ensemble": [0.05],
26 "selection_p_ensemble": [0.5],
27 "runcontrol_plotfitness": True,
28 "runcontrol_plotrankbest": True,
29 "runcontrol_plotrankall": False,
30 "runcontrol_plotsavefig": True,
31 }
32
33 gp = gpclassifier.initialize(**params)
34 gp.evolve()
35 gp.report(ensemble_row=0) # writes HTML report (e.g., ensemble_0.html)

```

## 1.4.2 Regression Task: gpregressor

**Listing 1.3:** Minimal regression example with gpregressor.

```

1 from genforge.gpregressor import gpregressor
2 import numpy as np
3
4 X_train, y_train = np.load("Xtr.npy"), np.load("ytr.npy")
5 X_val, y_val = np.load("Xva.npy"), np.load("yva.npy")
6 X_test, y_test = np.load("Xte.npy"), np.load("yte.npy")
7
8 params = {
9     "userdata_name": "demo_reg",
10    "userdata_xtrain": X_train, "userdata_ytrain": y_train,
11    "userdata_xval": X_val, "userdata_yval": y_val,
12    "userdata_xtest": X_test, "userdata_ytest": y_test,

```

```

13 # "userdata_pop_idx": findex, # (optional)
14 "runcontrol_pop_size": 64,
15 "runcontrol_generations": 50,
16 "runcontrol_useparallel": False,
17 "functions_name": [["times", "minus", "plus", "divide"]],
18 "gene_max_genes": [5],
19 "tree_max_depth": [10],
20 "selection_elite_fraction": [0.05],
21 "selection_elite_fraction_ensemble": [0.05],
22 "selection_p_ensemble": [0.5],
23 "runcontrol_plotfitness": True,
24 "runcontrol_plotrankbest": True,
25 "runcontrol_plotrankall": False,
26 "runcontrol_plotsavefig": True,
27 }
28
29 gp = gpregressor.initialize(**params)
30 gp.evolve()
31 gp.report(ensemble_row=0) # writes HTML report (e.g., ensemble_0.html)

```

### 1.4.3 Feature partitioning: `SPFPPartitioner`

**Listing 1.4:** Quickstart for SPFP feature partitioning.

```

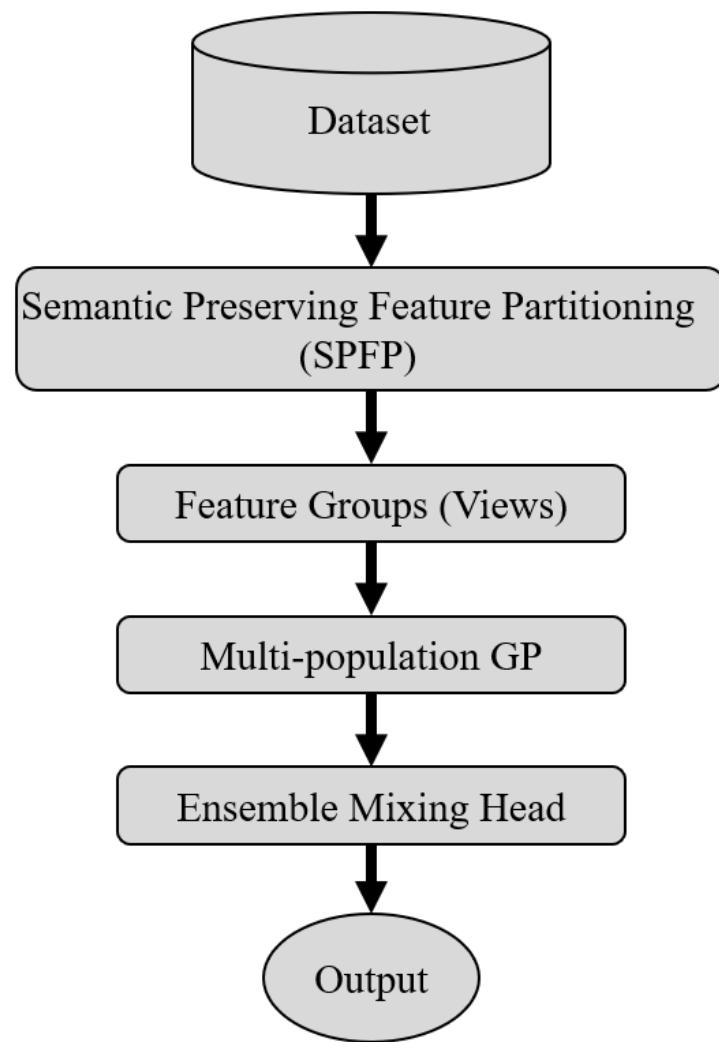
1 from genforge.spfp.spfp_partition import SPFPPartitioner
2 import numpy as np
3
4 X, y = np.load("X.npy"), np.load("y.npy")
5
6 spfp = SPFPPartitioner(
7     n_groups=5, # number of feature groups to construct
8     n_bins=100, # discretization bins
9     objective="micor", # "mi", "cor", or "micor"
10    override=False,
11    backward=False,
12    perfs=0.10, # minimum ratio of features per group
13    remp=0.60, # random removal ratio between groups
14    random_state=7,

```

```
15     verbose=1
16 )
17
18 spfp.fit(X, y) # optional: pass X_test, y_test for diagnostics
19 groups = spfp.partition() # build partitions
20 partitions = spfp.get_partitions() # [[fid...], [fid...], ...]
```

## 1.5 Reproducibility notes

Set `random_state` (where available) and/or any internal `SEED`, fix `num_threads`, and pin versions of dependencies. Genetic search is inherently stochastic; for apples-to-apples comparisons, keep all seeds and configuration identical across runs.



**Figure 1.1:** GenForge data flow: (1) optional SPFP partitions features into  $P$  groups ( $findex$ ); (2) each group maps to a GP population; (3) a linear mixing head combines population outputs into the final prediction.

## 2 Installation and Environment Setup

This chapter describes how to install `GenForge`, verify your environment, and prepare your workspace for running `gpclassifier`, `gpregressor`, and `SPFPPartitioner`.

### 2.1 System Requirements

**Table 2.1:** Supported environments and core dependencies.

Component	Recommended setup
Operating system	Linux, macOS, or Windows 10/11
Python	3.9–3.12
Core libraries	<code>numpy</code> , <code>scipy</code> , <code>joblib</code> , <code>matplotlib</code> (for plots)
Optional	<code>xgboost</code> (benchmarks), <code>pandas</code> (I/O convenience)
Not required	<code>scikit-learn</code> (replaced by internal <code>StandardScaler</code> , <code>LabelEncoder</code> )

**Build tools (source installs).** On Linux, ensure `gcc/g++` and `make`. On macOS, install Xcode Command Line Tools (`xcode-select -install`). On Windows, install Microsoft C++ Build Tools.

### 2.2 Install from PyPI

Currently `GenForge` is being considered to be published in PyPI. When `GenForge` is published on PyPI, install with:

**Listing 2.1:** Install from PyPI.

```
1 python -m pip install --upgrade pip
2 pip install genforge
3 # Optional extras (plots, xgboost benchmarks):
4 pip install matplotlib xgboost
```

## 2.3 Install from Source

Clone and install in editable mode:

**Listing 2.2:** Clone and editable install.

```

1 git clone https://github.com/maisamkhorshidi/genforge.git
2 cd genforge
3 python -m pip install --upgrade pip
4 pip install -e .[dev]
```

**Local development layout.** A typical tree:

**Listing 2.3:** Example project layout.

```

1 your-project/
2   data/
3     X_train.npy y_train.npy X_val.npy y_val.npy X_test.npy y_test.npy
4   runs/ # logs, reports, plots
5   scripts/ # experiment scripts
6   venv/ # optional virtual environment
```

## 2.4 Verifying the Installation

Run a quick import/version check:

**Listing 2.4:** Import and version check.

```

1 import genforge
2 print("GenForge version:", getattr(genforge, "__version__", "unknown"))
3
4 # Either top-level or submodule imports (depending on your __init__.py):
5 from genforge import gpclassifier, gpregressor, SPFPartitioner
6 # or:
7 # from genforge.gpclassifier import gpclassifier
8 # from genforge.gpregressor import gpregressor
9 # from genforge.spfp import SPFPartitioner
10
11 print("Imports OK.")
```

If imports succeed and the version prints, your environment is ready.

## 2.5 Parallelism and Reproducibility

**CPU parallelism.** Fitness evaluation in `GenForge` can run in parallel. Typical flags include `runcontrol_useparallel` and `runcontrol_n_jobs`. Set `runcontrol_batch_size` to balance work per worker.

**Determinism.** For reproducible results, set all available seeds: `random_state` in `SPFPPartitioner`, and any GP-level configuration seeds (e.g., for population initialization). Pin dependency versions in a `requirements.txt` or `pyproject.toml` for consistent environments.

## 2.6 Data Preparation

**Feature scaling.** Use the internal `StandardScaler` (`genforge.spfp`) to avoid external dependencies. For label encoding, use the internal `LabelEncoder`.

**Splits.** Provide NumPy arrays: `X_train` (N by D), `y_train` (N,), and analogous validation/test splits. If you plan to use `SPFPPartitioner`, run it on the same training split that your GP will consume.

## 2.7 Assets to Include for This Chapter

**Remark 2.1** (Asset checklist). Place the following under `Figure/` and tell us the filenames to reference them precisely.

1. Environment layout screenshot or diagram (e.g., `fig_install_layout.png`); we will reference it as Figure 2.1.
2. A short terminal capture of successful `pip install genforge` (optional; include as a figure or listing).

Name	Status	Date modified	Type	Size
__pycache__	✓	9/8/2025 7:14 AM	File folder	
gpclassifier	✓	9/8/2025 5:53 AM	File folder	
gpregressor	✓	9/8/2025 5:53 AM	File folder	
spfp	✓	9/8/2025 5:53 AM	File folder	
__init__.py	✓	9/8/2025 5:53 AM	Python Source File	1 KB

**Figure 2.1:** Example local project layout for GenForge experiments.

# 3 gpclassifier & gpregressor Configurations

## 3.1 Overview

gpclassifier and gpregressor configured through a typed container `ClassifierConfig` and `RegressorConfig`. Note that typed configurations are preferred over legacy dictionary for backward compatibility. At initialization the configuration is *validated and resolved*: defaults are applied, scalar values are broadcast to all populations where allowed, the function set is imported, and shapes/ranges are checked. The random seed in `runcontrol.random_state` is normalised into a NumPy `Generator` stored as `runcontrol.RNG`. Upon resolution during `gpclassifier(cfg)` and `gpregressor(cfg)` (equivalently, `gp.initialize(cfg)` when an initializer is exposed), the fully checked and expanded parameters are written to the runtime dictionaries `gp.config` and `gp.userdata` for the engine to consume; see Listing 3.1 and Listing 3.2 for a minimal import and initialization example.

The `ClassifierConfig` dataclass is organised into ten branches of parameters that are validated and then propagated to `gp.config / gp.userdata`: `runcontrol`, `user`, `selection`, `nodes`, `operator`, `gene`, `tree`, `fitness`, `softmax`, and `log` for classification task. Also, the `ClassifierConfig` dataclass is organised into ten branches of parameters that are validated and then propagated to `gp.config / gp.userdata`: `runcontrol`, `user`, `selection`, `nodes`, `operator`, `gene`, `tree`, `fitness`, `linregression`, and `log` for regression task. Note that many parameters and branches are common across `ClassifierConfig` and `RegressorConfig`. In this chapter, each branch is documented in turn, with the post-resolution (effective) types and semantics enforced by the resolver. This makes the behaviour of downstream components—population initialisation, genetic operators, heads, plotting and reporting—explicit and reproducible for any given configuration.

Unless noted otherwise, many parameters accept either a scalar or a boolean that is applied uniformly to every population, or a list whose length equals `runcontrol.num_pop`. During resolution, scalars and booleans are broadcast and inputs are coerced to the accepted forms listed below (for example, numeric

literals are cast to the target dtype, sequences are validated for length and range, and strings that name functions or optimisers are normalised against the supported sets). This broadcasting and coercion are part of the same resolution step that writes the final values into `gp.config` and `gp.userdata`, so what you see in those dictionaries is exactly what the training/evaluation code will use.

Unless noted otherwise, the parameters are common among `ClassifierConfig` and `RegressorConfig`.

**Listing 3.1:** Minimal usage of Importing `gpclassifier` and `ClassifierConfig` for classification task.

```

1 from genforge.gpclassifier import gpclassifier
2 from genforge.gpclassifier import ClassifierConfig
3
4 cfg = ClassifierConfig()
5
6 # Required user data (NumPy arrays)
7 cfg.user.xtrain = X_train # shape (N, D)
8 cfg.user.ytrain = y_train # shape (N,)
9 cfg.user.xval = X_val # optional, or None
10 cfg.user.yval = y_val # optional, or None
11 cfg.user.xtest = X_test # optional, or None
12 cfg.user.ytest = y_test # optional, or None
13
14 # Minimal run control
15 cfg.runcontrol.pop_size = 64
16 cfg.runcontrol.generations = 50
17 cfg.runcontrol.random_state = 7 # int / None / np.random.Generator / np.random.RandomState
18
19 gp = gpclassifier.initialize(cfg) # resolves and validates, and initializes
20 gp.evolve()
21 gp.report(ensemble_row=0) # writes an HTML report

```

## 3.2 Type conventions used for `gpclassifier`

To keep the notation compact:

- `int`, `float`, `bool`, `str` are Python types.
- `list[int]` means a Python list of `int`.

- “scalar or list[.]” means a single value broadcast to all populations or a per-population list whose length equals `runcontrol.num_pop`.
- Arrays are NumPy `ndarray`.

**Listing 3.2:** Minimal usage of Importing `gpregressor` and `RegressorConfig` for classification task.

```

1 from genforge.gpregressor import gpregressor
2 from genforge.gpregressor import RegressorConfig
3
4 cfg = RegressorConfig()
5
6 # Required user data (NumPy arrays)
7 cfg.user.xtrain = X_train # shape (N, D)
8 cfg.user.ytrain = y_train # shape (N,)
9 cfg.user.xval = X_val # optional, or None
10 cfg.user.yval = y_val # optional, or None
11 cfg.user.xtest = X_test # optional, or None
12 cfg.user.ytest = y_test # optional, or None
13
14 # Minimal run control
15 cfg.runcontrol.pop_size = 64
16 cfg.runcontrol.generations = 50
17 cfg.runcontrol.random_state = 7 # int / None / np.random.Generator / np.random.RandomState
18
19 gp = gpregressor.initialize(cfg) # resolves and validates, and initializes
20 gp.evolve()
21 gp.report(ensemble_row=0) # writes an HTML report

```

### 3.3 Type conventions used for `gpregressor`

To keep the notation compact:

- `int`, `float`, `bool`, `str` are Python types.
- `list[int]` means a Python list of `int`.
- “scalar or list[.]” means a single value broadcast to all populations or a per-population list whose length equals `runcontrol.num_pop`.

- Arrays are NumPy `ndarray`.

## 3.4 User Data

`user.name (str)`. (Optional, default: *Example\_Run*) A run label used in logs and reports.

`user.xtrain (ndarray 2D)`. (Required) Training features of shape `(N, D)`. Enforced to be exactly 2D.

`user.ytrain (ndarray 1D or 2D single column)`. (Required) Training labels. Internally flattened to 1D; labels are remapped to `0..K-1` if needed (a warning is issued).

`user.xval, user.yval (ndarray or None)`. (Optional) Validation split. Shapes are checked; if present, the number of columns must match `user.xtrain`.

`user.xtest, user.ytest (ndarray or None)`. (Optional) Test split with the same constraints.

`user.pop_idx (list[int] or list[list[int]] or None)`. (Optional, default: `[0,...,len(user.xtrain)]`) Feature indices per population. If multiple populations are used, this should be `list[list[int]]` with valid column indices per group.

`user.initial_population (list or None)`. (Optional, default: `None`) Initial population per population id; when provided the outer list must have length `runcontrol.num_pop`.

`user.stats (bool)`. (Optional, default: `True`) Collect additional runtime diagnostics for reporting.

## 3.5 Run Control

`runcontrol.num_pop (int)`. (Optional, default: `None`) Number of GP populations. If `user.pop_idx` is a list of groups, it resolves to its length.

`runcontrol.pop_size (int)`. (Optional, default: `25`) Individuals per population.

`runcontrol.generations (int)`. (Optional, default: `150`) Maximum generations to evolve.

`runcontrol.batch_job (int)`. (Optional, default: `5`) Batch size per worker in parallel mode.

**runcontrol.stallgen (int).** (Optional, default: 20) Early-stopping window (no improvement) used for diagnostics.

**runcontrol.verbose (int).** (Optional, default: 1) Console print frequency / verbosity.

**runcontrol.quiet (bool).** (Optional, default: False) Quiet logging (no console output).

**runcontrol.useparallel (bool).** (Optional, default: False) Enable parallel fitness evaluation.

**runcontrol.n\_jobs (int).** (Optional, default: 1) Worker processes if parallel is enabled.

**runcontrol.usecache (bool).** (Optional, default: True) Use fitness cache when cloning individuals.

**runcontrol.minimisation (bool).** (Optional, default: True) Minimise fitness (else maximise).

**runcontrol.tolfit (float).** (Optional, default: 1e-9) Improvement tolerance for convergence checks.

**runcontrol.plotfitness (bool).** (Optional, default: False) Toggle Fitness vs. Generation plot.

**runcontrol.plotrankall (bool).** (Optional, default: False) Toggle RankAll plots.

**runcontrol.plotrankbest (bool).** (Optional, default: False) Toggle RankBest plots.

**runcontrol.plotformat (str or list[str]).** (Optional, default: "png") Image format(s) to save.

**runcontrol.plotfolder (str or None).** (Optional, default: None) Plot output folder; if None resolves to CWD.

**runcontrol.plotsavefig (bool).** (Optional, default: True) Save figures to disk.

**runcontrol.plotlive (bool).** (Optional, default: False) Live drawing during the run (if supported).

**runcontrol.plotbackend (str).** (Optional, default: "auto") Matplotlib backend: e.g., "Agg", "QtAgg", "TkAgg", "MacOSX".

**runcontrol.track\_individuals (bool).** (Optional, default: False) Track individuals across generations (memory-intensive).

**runcontrol.resultfolder (str or None).** (Optional, default: `None`) Output folder for run artifacts; if `None` resolves to CWD.

**run.adaptgen (bool).** (Specific to `RegressorConfig`) (Optional, default: `False`) Activates the adaptive hierarchical symbolic abstraction (AHSAM) module.

**runcontrol.adaptgen (int).** (Specific to `RegressorConfig`) (Optional, default: `15`) The number of generations at which, after no fitness improvement, the AHSAM is performed.

### Randomness and reproducibility.

**runcontrol.random\_state (None or int or np.random.Generator or np.random.RandomState).** (Optional, default: `None`) User-provided seed/state kept for provenance.

**runcontrol.RNG (np.random.Generator).** (Resolved, no user default) Created at initialisation from `runcontrol.random_state`; the engine routes all randomness through this generator.

## 3.6 Selection

(All entries accept a scalar or a per-population list unless noted.)

**selection.tournament\_size (int).** (Optional, default: `2`) Tournament size.

**selection.elite\_fraction (float).** (Optional, default: `0.05`) Elite fraction within a population.

**selection.elite\_fraction\_ensemble (float).** (Optional, default: `0.05`) Elite fraction for cross-population ensemble building.

**selection.tournament\_lex\_pressure (bool).** (Optional, default: `True`) Lexicographic tie-breaking pressure.

**selection.tournament\_p\_pareto (float).** (Optional, default: `0.0`) Probability of using a Pareto tournament.

**selection.p\_ensemble (float).** (Optional, default: `0.0`) Probability of sampling ensemble fitness during selection.

## 3.7 Nodes: Ephemeral Constants and Function Set

### Ephemeral Random Constants (ERCs)

**nodes.const.about (str).** (Optional, default: "Ephemeral random constants") Informational string (preserved).

**nodes.const.p\_ERC (float or list).** (Optional, default: 0.1) Probability of generating an ERC when creating a leaf.

**nodes.const.p\_int (float or list).** (Optional, default: 0.0) Probability that an ERC is integer-typed.

**nodes.const.range (list[2] or list of pairs).** (Optional, default: [-10, 10]) Sampling range; each [min, max] must satisfy min < max.

**nodes.const.num\_dec\_places (int or list).** (Optional, default: 4) Decimal places for rounding real-valued ERCs.

### Function Set

**nodes.functions.name (list[str] or list[list[str]]).** (Required) Operator names per population. Should be provided by the user as the name of corresponding operator script in the working path, e.g. ["times", "minus", "plus"]

## 3.8 Genetic Operators

**operator.p\_mutate (float or list).** (Optional, default: 0.14) Probability of mutation.

**operator.p\_cross (float or list).** (Optional, default: 0.84) Probability of crossover.

**operator.p\_direct (float or list).** (Optional, default: 0.02) Probability of direct reproduction.

**operator.mutate\_par (list[float] of length 6, or per-pop list).** (Optional, default: [0.9, 0.05, 0.05, 0.0, 0.0, 0.0]) Sub-operator probabilities (normalised if sum ≠ 1).

**operator.mutate\_par\_cumsum (list[list[float]]).** (Resolved, created at initialisation) Cumulative sums used for fast sampling.

**operator.mutate\_gaussian\_std (float or list).** (Optional, default: 0.1) Std. dev. for Gaussian perturbation of numeric ERCs.

## 3.9 Gene-Level Controls

**gene.p\_cross\_hi (float or list).** (Optional, default: 0.2) Probability of high-level (inter-gene) crossover.

**gene.hi\_cross\_rate (float or list).** (Optional, default: 0.5) Probability any given gene participates in high-level crossover.

**gene.multigene (bool or list).** (Optional, default: True) Allow multiple genes per individual.

**gene.max\_genes (int or list).** (Optional, default: 5) Maximum number of genes per individual.

## 3.10 Tree Constraints

**tree.build\_method (int or list).** (Optional, default: 3) Tree builder method id 3 = ramped half-and-half, 2 = grow, 1 = full.

**tree.max\_nodes (int or np.inf or list).** (Optional, default: np.Inf) Maximum node count per tree.

**tree.max\_depth (int or list).** (Optional, default: 4) Maximum tree depth.

**tree.max\_mutate\_depth (int or list).** (Optional, default: 4) Depth cap for mutation-generated subtrees.

## 3.11 Fitness and Objectives

**fitness.terminate (bool).** (Optional, default: False) Enable early termination when threshold met.

**fitness.complexityMeasure (int).** (Optional, default: 1) Complexity penalty mode: 1 (expressional complexity) or 0 (node count).

## 3.12 Softmax Head (Specific to `ClassifierConfig`)

These parameters control the softmax head used in `gpclassifier`. Parameters are broadcast per population unless a list is provided.

**softmax.optimizer\_type** (`str or list`). (Optional, default: `rmsprop`) Name of the optimizer. Allowed values include `sgd`, `adam`, `rmsprop`, `adamw`, `nadam`, `adagrad`, `adadelta`, `ftrl`, `adamax`, `sgdnm`, `lbfgs`.

**softmax.optimizer\_param** (`None or dict or list[dict]`). (Optional, default: `None`) Per-optimizer hyperparameters merged onto validated defaults. Unknown keys are ignored with a warning. Examples: `learning_rate`, `beta1`, `beta2`, `epsilon`, `weight_decay`.

**softmax.initializer** (`str or list`). (Optional, default: `glorot_uniform`) Weight/bias initializer family (e.g., `glorot_uniform`, `he_normal`, `random_normal`).

**softmax.regularization** (`None or str or list`). (Optional, default: `None`) Regularization type: `l1`, `l2`, or `hybrid`. If `hybrid` is used with a scalar `regularization_rate`, a warning is issued (`hybrid` expects a pair).

**softmax.regularization\_rate** (`float or list`). (Optional, default: `0.01`) Regularization strength. For `softmax.regularization=hybrid`, pass a pair like `[l1, l2]`.

**softmax.batch\_size** (`int or list`). (Optional, default: `32`) Mini-batch size.

**softmax.epochs** (`int or list`). (Optional, default: `1000`) Maximum training epochs.

**softmax.patience** (`int or list`). (Optional, default: `10`) Early-stopping patience.

**softmax.buffer\_size** (`None or int or list`). (Optional, default: `None`) Optional shuffle/prefetch buffer size.

**softmax.shuffle** (`bool or list`). (Optional, default: `True`) Shuffle batches each epoch.

**softmax.verbose** (`int or list`). (Optional, default: `0`) Trainer verbosity.

### 3.13 ElasticNetCV Head (Specific to `RegressorConfig`)

These parameters control the ElasticNetCV regression head used in `gpregressor`. Parameters are broadcast per population unless a list is provided.

**`linregression.l1_ratio` (float or list).** (Optional, default: 1.0) Mixing ratio between L1 and L2 penalties: 0.0 corresponds to Ridge, 1.0 to Lasso, intermediate values interpolate to Elastic Net.

**`linregression.alphas` (None or list[float] or list of lists).** (Optional, default: None) Explicit list of alpha values to try. If provided, this overrides `linregression.eps` and `linregression.n_alphas`.

**`linregression.n_alphas` (int or list).** (Optional, default: 100) Number of alpha values in the automatically generated grid (ignored if `linregression.alphas` is set).

**`linregression.eps` (float or list).** (Optional, default: 1e-3) When `linregression.alphas` is None, the alpha grid is logarithmically spaced between  $\alpha_{\min} = \text{eps} \times \alpha_{\max}$  and  $\alpha_{\max}$ .

**`linregression.fit_intercept` (bool or list).** (Optional, default: True) Whether to fit an intercept. If features and targets are already centered, this can be set to False.

**`linregression.copy_x` (bool or list).** (Optional, default: True) If True, training data is copied; if False, it may be overwritten/centered in place (small performance gain, use with care).

**`linregression.max_iter` (int or list).** (Optional, default: 1000) Maximum coordinate-descent iterations per fit. Increase if convergence warnings appear.

**`linregression.tol` (float or list).** (Optional, default: 1e-4) Optimization stopping tolerance on the objective. Smaller values yield more precise but slower convergence.

**`linregression.cv` (int or list).** (Optional, default: 5) Number of cross-validation folds.

**`linregression.n_jobs` (None or int or list).** (Optional, default: None) Number of parallel jobs for cross-validation. None means 1 unless in `joblib.parallel_backend` context.

**`linregression.verbose` (int or list).** (Optional, default: 0) Verbosity level. 0 = silent, higher values increase output.

**linregression.positive** (`bool or list`). (Optional, default: `False`) Constrain coefficients to be non-negative. Useful for interpretability but can worsen fit if data are not scaled/centered.

**linregression.selection** (`str or list`). (Optional, default: `"cyclic"`) Coordinate-descent sweep order. Allowed values:

- `"cyclic"` (deterministic sweeps),
- `"random"` (randomised order, sometimes faster on large problems).

**Head seed handling.** The softmax head derives its reproducible integer seed from the resolved `runcontrol.RNG`; it does not accept a separate external seed. If `runcontrol.random_state` is a NumPy `Generator` (or a legacy `RandomState`), the integer is drawn from that source so that repeated runs with the same RNG are reproducible. If `runcontrol.random_state` is `None` or an `int`, behavior follows standard NumPy semantics.

## 3.14 Resolved dictionaries

After validation, two internal dictionaries are constructed:

**config** (`dict`). Resolved engine configuration with sections `runcontrol`, `selection`, `nodes`, `operator`, `gene`, `tree`, `fitness`, and `softmax`. Values are fully broadcast and type-checked.

**userdata** (`dict`). Immutable copies of user arrays and metadata: `name`, `stats`, population indices, `xtrain/xval/xtest`, `ytrain/yval/ytest`, and binarised labels used by the classifier.

**Label remapping.** If the unique sorted labels in `user.ytrain` are not `0..K-1`, they are remapped internally (a warning is emitted). Predictions and reports use the remapped indexing unless you invert-map externally.

## 3.15 Common validation rules (summary)

- Arrays must have expected dimensionality: features strictly 2D; labels 1D or a single 2D column (flattened internally).

- Probabilities must lie in `[0, 1]` unless specified otherwise.
- ERC ranges `[min, max]` must satisfy `min < max`.
- Lists provided for per-population parameters must match `runcontrol.num_pop`; scalars are broadcast.
- Unknown optimizer types are rejected; unknown optimizer hyperparameters are warned and ignored.
- No global `np.random` seeding: all randomness flows through `runcontrol.RNG`.

## 3.16 Complete example

**Listing 3.3:** Putting it together: a compact configuration for two populations.

```

1 from genforge.gpclassifier import gpclassifier, ClassifierConfig
2 import numpy as np
3
4 cfg = ClassifierConfig()
5
6 # data
7 cfg.user.name = "demo_clf_two_pop"
8 cfg.user.xtrain = X_train
9 cfg.user.ytrain = y_train
10 cfg.user.xval = X_val
11 cfg.user.yval = y_val
12 cfg.user.xtest = X_test
13 cfg.user.ytest = y_test
14
15 # tie features to populations (two groups)
16 cfg.user.pop_idx = [list(range(0, 20)), list(range(20, 40))]
17
18 # run control
19 cfg.runcontrol.num_pop = 2
20 cfg.runcontrol.pop_size = 64
21 cfg.runcontrol.generations = 50
22 cfg.runcontrol.useparallel = True

```

```

23 cfg.runcontrol.parallel.n_jobs = 4
24 cfg.runcontrol.parallel.batch_job = 8
25 cfg.runcontrol.random_state = 123
26
27 # nodes/functions
28 cfg.nodes.functions.name = [
29     ["times", "minus", "plus", "divide"],
30     ["times", "minus", "plus", "divide", "sin"]
31 ]
32 cfg.nodes.const.p_ERC = [0.3, 0.3]
33 cfg.nodes.const.p_int = [0.2, 0.2]
34 cfg.nodes.const.range = [[-3.0, 3.0], [-1.0, 1.0]]
35 cfg.nodes.const.num_dec_places = [3, 3]
36
37 # selection
38 cfg.selection.tournament_size = [5, 5]
39 cfg.selection.elite_fraction = [0.05, 0.05]
40 cfg.selection.elite_fraction_ensemble = [0.05, 0.05]
41 cfg.selection.tournament_lex_pressure = [True, True]
42 cfg.selection.tournament_p_pareto = [0.5, 0.5]
43 cfg.selection.p_ensemble = [0.5, 0.5]
44
45 # operator
46 cfg.operator.p_mutate = [0.40, 0.40]
47 cfg.operator.p_cross = [0.50, 0.50]
48 cfg.operator.p_direct = [0.10, 0.10]
49 cfg.operator.mutate_par = [[0.20, 0.20, 0.20, 0.20, 0.10, 0.10]] # broadcasted
50 cfg.operator.mutate_gaussian_std = [0.10, 0.10]
51
52 # gene / tree
53 cfg.gene.multigene = [True, True]
54 cfg.gene.max_genes = [5, 5]
55 cfg.gene.p_cross_hi = [0.3, 0.3]
56 cfg.gene.hi_cross_rate = [0.5, 0.5]
57
58 cfg.tree.build_method = [0, 0]
59 cfg.tree.max_nodes = [np.inf, np.inf]
60 cfg.tree.max_depth = [10, 10]
```

```
61 cfg.tree.max_mutate_depth = [8, 8]
62
63 # fitness
64 cfg.fitness.terminate = False
65 cfg.fitness.complexityMeasure = 0
66
67 # softmax head
68 cfg.softmax.optimizer_type = ["adam", "adam"]
69 cfg.softmax.optimizer_param = [{"learning_rate": 0.001}, {"learning_rate": 0.001}]
70 cfg.softmax.initializer = ["glorot_uniform", "glorot_uniform"]
71 cfg.softmax.regularization = ["l2", "l2"]
72 cfg.softmax.regularization_rate = [1e-4, 1e-4]
73 cfg.softmax.batch_size = [64, 64]
74 cfg.softmax.epochs = [100, 100]
75 cfg.softmax.patience = [10, 10]
76 cfg.softmax.buffer_size = [None, None]
77 cfg.softmax.shuffle = [True, True]
78 cfg.softmax.verbose = [0, 0]
79
80
81 gp = gpclassifier.initialize(cfg)
82 gp.evolve()
```

## 4 SPFP: Mathematics and Algorithmic Details

Semantic Preserving Feature Partitioning (`SPFPPartitioner`) constructs  $P$  feature groups (`findex`) to drive multi-population genetic programming. This chapter formalises the notation, discretisation, scoring, forward selection with growth constraints, optional backward pruning, and the final group formation policy.

### 4.1 Notation and Discretisation

Let  $X \in \mathbb{R}^{N \times D}$  be the input matrix and  $y \in \mathbb{R}^N$  the target. The pipeline standardises  $X$  feature-wise to zero mean and unit variance, and label-encodes  $y$  if needed. For entropy-based quantities, variables are discretised into  $B$  equal-distance bins (`n_bins`).

For any index set  $S \subseteq \{1, \dots, D\}$ , define the discrete random vector  $X_S$  as the concatenation of the binned features  $\{X_j\}_{j \in S}$ . Denote the joint entropy

$$H(X_S) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}), \quad (4.1)$$

and similarly  $H(X_S, y)$  and  $H(y)$ . `SPFPPartitioner` uses the normalised mutual information (NMI) with respect to  $H(y)$ :

$$\text{NMI}(X_S; y) = \frac{H(X_S) + H(y) - H(X_S, y)}{H(y)}. \quad (4.2)$$

All joint entropies are computed by a joint-histogram routine with consistent bin edges (training edges reused for test, when given).

## 4.2 Per-Feature Scoring and Ranking

For each feature  $i \in \{1, \dots, D\}$ , compute

$$H_i := H(X_{\{i\}}), \quad (4.3)$$

$$H_{i,y} := H(X_{\{i\}}, y), \quad (4.4)$$

$$\text{mi}_i := \frac{H_i + H(y) - H_{i,y}}{H(y)}. \quad (4.5)$$

On the *continuous, standardised* features and label-encoded  $y$ , compute the Pearson correlation  $\rho_i$ , and define the auxiliary score

$$\text{micor}_i = \text{mi}_i + |\rho_i|. \quad (4.6)$$

The initial *stream* of candidate features is sorted by  $\text{mi}$  in descending order.<sup>1</sup>

## 4.3 Group-Wise Forward Selection

SPFPPartitioner builds  $P = \text{num\_groups}$  groups sequentially. At the beginning of a group, let  $A$  be the set of remaining candidates (a copy of the current stream). Define the group-level reference values

$$H_A := H(X_A), \quad H_{A,y} := H(X_A, y), \quad \text{NMI}_A := \frac{H_A + H(y) - H_{A,y}}{H(y)}. \quad (4.7)$$

Let  $S$  be the set of *selected* indices for the current group (initially empty). Define the running quantities  $H_S := H(X_S)$  and  $H_{S,y} := H(X_S, y)$ , and the current normalised mutual information

$$\text{NMI}_S = \frac{H_S + H(y) - H_{S,y}}{H(y)}. \quad (4.8)$$

**Stopping rule.** Selection continues while

$$\left( H_{S,y} < \alpha \cdot \lfloor H_{A,y} \rfloor_\delta \right) \vee \left( H_S < \alpha \cdot \lfloor H_A \rfloor_\delta \right) \vee \left( |S| < \lceil \pi \cdot D \rceil \right), \quad (4.9)$$

and there remain candidates in  $A$ . Here  $\alpha = \text{criterion\_add}(\text{wcriadd})$ ,  $\pi = \text{perfs}$ , and  $\lfloor \cdot \rfloor_\delta$  is a fixed-decimal floor with  $\delta = 10^{-\text{floor\_decimals}}$  (the reference uses 4 decimals in this context).

<sup>1</sup>This mirrors the reference script, which sorts by  $\text{mi}$  and then constructs groups from that stream.

**First pick.** On the first iteration, select the top-ranked candidate in  $A$  (the head of the stream), add it to  $S$ , and update  $H_S, H_{S,y}$ .

**Composite objective.** From the second iteration onward, maintain per-candidate components for each  $j \in A$  (with current  $S \neq \emptyset$ ):

$$\text{objmi}_j := \text{mi}_j, \quad \text{objcov}_j := |\rho_j|, \quad (4.10)$$

$$\text{objint}_j := \frac{1}{|S|} \sum_{k \in S} \frac{-H(y) - H(X_{\{j,k\}}, y) + H_{j,y} + H_{k,y}}{H(y)}, \quad (4.11)$$

$$\text{objred}_j := \frac{1}{|S|} \sum_{k \in S} \frac{H_j + H_k - H(X_{\{j,k\}})}{H(y)}, \quad (4.12)$$

and the additive composite

$$\text{obj}_j = \text{objmi}_j + \text{objcov}_j + \text{objint}_j + \text{objred}_j. \quad (4.13)$$

Equations (4.11) and (4.12) mirror the reference implementation: conditional interaction uses  $I(X_j; X_k | y)$  scaled by  $H(y)$ , and redundancy uses  $I(X_j; X_k)$  scaled by  $H(y)$ ; both are averaged over the currently selected set  $S$ .

**Selection policy.** On iteration  $t \geq 2$ , `SPFPPartitioner` chooses the next candidate by one of the following branches (exactly as implemented):

- (i) **Objective gate.** If both coverage thresholds in (4.9) are already met but  $|S| < \lceil \pi D \rceil$ , choose the maximiser of a user-chosen primary score `objective`  $\in \{ \text{"mi"}, \text{"cor"}, \text{"micor"} \}$ :

$$j^* \in \arg \max_{j \in A} \begin{cases} \text{mi}_j, & \text{"mi"} \\ |\rho_j|, & \text{"cor"} \\ \text{mi}_j + |\rho_j|, & \text{"micor"} \end{cases}$$

- (ii) **Override.** If `override` is `True`, choose  $j^* \in \arg \max_{j \in A} \text{obj}_j$ .

- (iii) **Growth-gated search.** Otherwise, scan candidates in decreasing  $\text{obj}_j$  order while enforcing

conservative growth:

$$d_S(j) := \frac{H(X_{S \cup \{j\}}) - H_S}{\alpha \cdot \lfloor H_A \rfloor_\delta}, \quad d_{S,y}(j) := \frac{H(X_{S \cup \{j\}}, y) - H_{S,y}}{\alpha \cdot \lfloor H_{A,y} \rfloor_\delta}. \quad (4.14)$$

For the running maxima  $H_S$  and  $H_{S,y}$ , let  $\gamma = \text{growth_magnitude (whmag)}$ . If

$$\left( H_{S,y} < \alpha \lfloor H_{A,y} \rfloor_\delta \text{ or } H_S < \alpha \lfloor H_A \rfloor_\delta \right) \text{ and } \left( H(X_{S \cup \{j\}}, y) \leq \min(\gamma H_{S,y}, \alpha \lfloor H_{A,y} \rfloor_\delta) \text{ or } H(X_{S \cup \{j\}}) \leq \min(H_S, \alpha \lfloor H_A \rfloor_\delta) \right)$$

record  $d_S(j)$  and  $d_{S,y}(j)$  and *temporarily* suppress  $\text{obj}_j$  (set to  $-\infty$ ) to continue scanning. After the scan:

$$j^* = \begin{cases} \arg \max_{j: d_S(j) > 0 \wedge d_{S,y}(j) > 0} (d_S(j) + d_{S,y}(j)), & \text{if any} \\ \arg \max_{j: d_{S,y}(j) > 0} d_{S,y}(j), & \text{else if any} \\ \arg \max_{j: d_S(j) > 0} d_S(j), & \text{else if any} \\ \text{first scanned candidate,} & \text{otherwise.} \end{cases}$$

Once  $j^*$  is chosen, update  $S \leftarrow S \cup \{j^*\}$ , recompute  $H_S$ ,  $H_{S,y}$ , and refresh (4.11)–(4.13). This loop continues until the stopping rule (4.9) fails.

## 4.4 Optional Backward Pruning

If enabled (`backward = True`), pruning iteratively removes features from  $S$  as long as:

$$H_{S,y} \geq \beta \cdot \lfloor H_{S,y}^{\text{base}} \rfloor_\delta, \quad H_S \geq \beta \cdot \lfloor H_S^{\text{base}} \rfloor_\delta, \quad |S| > \lceil \pi D \rceil, \quad (4.15)$$

where  $\beta = \text{criterion_rem (wcrirem)}$ , and  $H_S^{\text{base}}, H_{S,y}^{\text{base}}$  are computed on the full current  $S$  before pruning starts.

For each  $k \in S$ , form the leave-one-out set  $S \setminus \{k\}$ , compute its entropies (consistent with forward definitions), its normalised mutual information, and interaction/redundancy summaries analogous to (4.11)–(4.12). Let

$$\text{obj}_k^{\text{back}} = \text{NMI}_{S \setminus \{k\}} + \frac{1}{|S| - 1} \sum I(\cdot; \cdot \mid y) - \frac{1}{|S| - 1} \sum I(\cdot; \cdot), \quad (4.16)$$

mirroring the script's sign convention in pruning.<sup>2</sup> Define admissible removals via threshold tests on  $H_S$  and  $H_{S,y}$  (and an additional NMI floor tied to  $\text{NMI}_A$  in (4.7)), and remove the  $k$  with the *smallest*  $\text{obj}_k^{\text{back}}$  among admissible candidates. Repeat while (4.15) holds.

## 4.5 Forming Multiple Groups and Partition Policy

At the end of a group, `SPFPPartitioner` samples without replacement a fraction  $\rho = \text{remp}$  of the selected indices  $S$  and removes those from the global stream. The remaining features stay available for subsequent groups. This random withholding encourages disjointness between groups and spreads high-value features across populations.

After  $P$  groups or exhaustion of the stream, `SPFPPartitioner` returns

$$\mathcal{P} = \{S_1, S_2, \dots, S_P\},$$

where each  $S_g$  is the (possibly empty) set of feature indices selected for group  $g$ .

## 4.6 Parameter Summary

**Table 4.1:** SPFP key parameters and recommended defaults (matching the reference implementation).

Name	Meaning	Typical default
<code>num_groups</code>	Number of groups $P$ to construct	5
<code>n_bins</code>	Bins for equal-distance discretisation	100
<code>objective</code>	Primary score in branch (i)	"mi" or "micor"
<code>criterion_add</code>	$\alpha$ in (4.9)	0.999–1.0
<code>criterion_rem</code>	$\beta$ in (4.15)	0.999–1.0
<code>growth_magnitude</code>	$\gamma$ growth cap in gated search	1.01
<code>perfs</code>	$\pi$ minimal ratio of $D$ per group	0.10 (dataset-dependent)
<code>remp</code>	$\rho$ random withholding rate	0.60
<code>override</code>	Use composite objective directly	<code>False</code>
<code>backward</code>	Enable pruning phase	<code>False</code>
<code>random_state</code>	Seed for reproducibility	(set as needed)

<sup>2</sup>In the reference script, the redundancy term is subtracted during backward pruning and added during forward selection.

## 4.7 Complexity and Implementation Notes

- **Entropy evaluations.** The dominant cost arises from repeated joint entropies such as  $H(X_{S \cup \{j\}})$ ,  $H(X_{S \cup \{j\}}, y)$ ,  $H(X_{\{j,k\}})$ , and  $H(X_{\{j,k\}}, y)$ . The reference uses parallel map for the pairwise terms involving the most recent selection.
- **Normalisation by  $H(y)$ .** All MI-derived terms are scaled by  $H(y)$ , exactly as in the script, to keep units comparable across datasets.
- **Flooring.** Threshold comparisons use decimal flooring on  $H_A$  and  $H_{A,y}$  with 4 decimals (see (4.9)).
- **Correlation.** Correlations are computed on the continuous, standardised  $X$  against label-encoded  $y$ ; only  $|\rho|$  is used.
- **Forward vs. backward signs.** Redundancy contributes positively in forward selection (added to (4.13)) but is penalised in pruning (subtracted in (4.16)).

## 4.8 Pseudocode

**Listing 4.1:** SPFP high-level pseudocode (mirrors the reference).

```

1 # Inputs: X (N x D), y (N,), n_groups P, n_bins, objective, criterion_add alpha,
2 # criterion_rem beta, growth_magnitude gamma, perfs pi, remp rho, override flag,
3 # backward flag, random_state.
4 # Pre: standardise X; label-encode y; bin X, y for entropy ops.
5
6 # 1) Per-feature scores
7 for i in 1..D:
8     Hx[i] = H({i})
9     Hxy[i] = H({i} + {y})
10    mi[i] = (Hx[i] + Hy - Hxy[i]) / Hy
11    cor[i] = abs(pearson(x_i, y_cont))
12    micor[i] = mi[i] + cor[i]
13    stream = argsort(mi)[desc] # initial order
14
15 # 2) Build P groups
16 for g in 1..P:
17     A = list(stream) # remaining candidates

```

```

18     S = [] # selected indices
19     H_A = H(A); H_Ay = H(A + {y})
20     while (H(S,y) < alpha*floor(H_Ay) or H(S) < alpha*floor(H_A) or len(S) < ceil(pi*D)) and A:
21         if len(S) == 0:
22             j_star = A[0]
23         else:
24             update objective components objmi, objcor, objint, objred
25             obj = objmi + objcor + objint + objred
26             if thresholds met but len(S) < ceil(pi*D):
27                 j_star = argmax objective_score in {"mi", "cor", "micor"}
28             elif override:
29                 j_star = argmax obj
30             else:
31                 gated-scan over sorted obj to compute dS, dSy per candidate
32                 choose by (dS>0 and dSy>0) -> max dS+dSy
33                 elif dSy>0 -> max dSy
34                 elif dS>0 -> max dS
35                 else fallback to first scanned
36             S.append(j_star); A.remove(j_star)
37             refresh H(S), H(S,y), and pairwise {H(j_star,f), H(j_star,f,y)} for f in A
38             update interaction/redundancy matrices and objective components
39
40     # Optional backward pruning
41     if backward:
42         while guards with beta thresholds and len(S) > ceil(pi*D):
43             compute leave-one-out metrics and back-objective
44             remove admissible k with smallest back-objective
45
46     # Partition policy: randomly remove floor(rho*len(S)) from global stream
47     remove_sample = random.sample(S, floor(rho*len(S)))
48     stream = stream \ remove_sample
49     save group S
50     return {S_1, ..., S_P}

```



# 5 gpclassifier: Legacy Configuration & API

This chapter documents the `gpclassifier` module: data requirements, the training pipeline, configuration keys (with defaults), integration with `SPFPPartitioner`, and the expected outputs (plots and HTML reports). *Note that typed configurations are preferred over legacy dictionary for backward compatibility*

## 5.1 Overview & Data Requirements for `gpclassifier`

**Task.** `gpclassifier` implements multi-population Genetic Programming (GP) for supervised classification. Multiple co-evolving populations contribute to an ensemble head that mixes population outputs into final class scores.

**Inputs.** Provide NumPy arrays:

- `userdata_xtrain`: shape  $(N_{tr}, D)$ ,
- `userdata_ytrain`: shape  $(N_{tr},)$  with integer labels,
- (optional) `userdata_xval`, `userdata_yval`,
- (optional) `userdata_xtest`, `userdata_ytest`.

If you use feature partitions from `SPFPPartitioner`, pass them via `userdata_pop_idx` (see Section 5.4).

## 5.2 Training Pipeline

At a high level, a training call:

1. builds initial populations (optionally restricted by `userdata_pop_idx`),
2. evolves for `runcontrol_generations` using selection, crossover, and mutation under structural limits,

3. fits an ensemble mixing head on population outputs (using the validation split when available),
4. evaluates and emits artifacts (plots, HTML reports, run snapshots).

## 5.3 Configuration Keys

The initializer expects a flat dictionary of keys. Below we group the most commonly used keys by purpose. Values shown as *defaults* are representative and may be dataset-dependent.

### 5.3.1 User data & run metadata

**Table 5.1:** User data and metadata keys for `gpclassifier`.

Key	Default	Description
<code>userdata_name</code>	"run"	Logical run name; used in filenames.
<code>userdata_xtrain</code>	—	Training features, ( $N_{tr}, D$ ).
<code>userdata_ytrain</code>	—	Training labels, ( $N_{tr}, .$ ).
<code>userdata_xval</code>	None	Validation features (optional).
<code>userdata_yval</code>	None	Validation labels (optional).
<code>userdata_xtest</code>	None	Test features (optional).
<code>userdata_ytest</code>	None	Test labels (optional).
<code>userdata_pop_idx</code>	None	Feature partitions (list of lists) from <code>SPFPPartitioner</code> .

### 5.3.2 Run control & parallelism

**Table 5.2:** Run control and parallel settings.

Key	Default	Description
<code>runcontrol_generations</code>	50	Number of generations.
<code>runcontrol_pop_size</code>	64	Individuals per population.
<code>runcontrol_useparallel</code>	False	Enable parallel fitness.
<code>runcontrol_n_jobs</code>	None	Workers (if parallel).
<code>runcontrol_batch_size</code>	None	Work chunk per worker.

### 5.3.3 Function sets, tree/gene constraints

**Table 5.3:** Structural configuration (function sets, depth, genes).

Key	Default	Description
functions_name	[["times", "minus", "plus", "divide"]]	Function set per population.
tree_max_depth	[10]	Max depth per population.
gene_max_genes	[5]	Max genes per chromosome (per population).

### 5.3.4 Selection & ensemble

**Table 5.4:** Selection and ensemble hyperparameters.

Key	Default	Description
selection_elite_fraction	[0.05]	Elitism per population.
selection_elite_fraction_ensemble	[0.05]	Elitism for ensemble stage.
selection_p_ensemble	[0.5]	Probability of choosing ensemble candidates.

### 5.3.5 Plotting & reporting

**Table 5.5:** Plot and report toggles.

Key	Default	Description
runcontrol_plotfitness	True	Plot Fitness vs. Generation.
runcontrol_plotrankbest	True	Plot RankBest curve.
runcontrol_plotrankall	False	Plot RankAll curves.
runcontrol_plotsavefig	True	Save plots to disk.

**Remark 5.1** (About defaults). Exact defaults can vary across releases. The keys above mirror the common usage patterns in the examples. When in doubt, consult the repository `README.md` and the source of `gp_class.py`.

## 5.4 Using SPFP Partitions

If you precompute feature partitions with `SPFPPartitioner` (see Chapter 4), pass them as a Python list of lists via `userdata_pop_idx`. The number of populations should match the number of groups:

**Listing 5.1:** Wiring SPFP partitions to populations.

```

1 from genforge.gpclassifier import gpclassifier
2 from genforge.spfp.spfp_partition import SPFPPartitioner
3
4 # X, y are NumPy arrays
5 spfp = SPFPPartitioner(n_groups=4, objective="micor", random_state=13)
6 spfp.fit(X, y)
7 findex = spfp.get_partitions() # e.g., [[0,7,9],[1,2,5],...]
8
9 params = {
10     "userdata_name": "clf_with_spfp",
11     "userdata_xtrain": X_train, "userdata_ytrain": y_train,
12     "userdata_xval": X_val, "userdata_yval": y_val,
13     "userdata_xtest": X_test, "userdata_ytest": y_test,
14     "userdata_pop_idx": findex, # <- map groups to populations
15     "runcontrol_generations": 50,
16     "runcontrol_pop_size": 64,
17     "functions_name": [["times", "minus", "plus", "divide"]],
18     "gene_max_genes": [5], "tree_max_depth": [10],
19     "selection_elite_fraction": [0.05],
20     "selection_elite_fraction_ensemble": [0.05],
21     "selection_p_ensemble": [0.5],
22     "runcontrol_plotfitness": True,
23     "runcontrol_plotrankbest": True,
24     "runcontrol_plotsavefig": True,
25 }
26 gp = gpclassifier.initialize(**params)
27 gp.evolve()
28 gp.report(ensemble_row=0)
```

## 5.5 Training & Evaluation Workflow

A minimal end-to-end script looks like Listing 5.2. Replace the paths to your arrays, and adjust `functions_name`, tree depth, and generations to suit your problem scale.

**Listing 5.2:** Minimal `gpclassifier` workflow.

```

1 from genforge.gpclassifier import gpclassifier
2 import numpy as np
3
4 X_train, y_train = np.load("X_train.npy"), np.load("y_train.npy")
5 X_val, y_val = np.load("X_val.npy"), np.load("y_val.npy")
6 X_test, y_test = np.load("X_test.npy"), np.load("y_test.npy")
7
8 params = {
9     "userdata_name": "demo_clf",
10    "userdata_xtrain": X_train, "userdata_ytrain": y_train,
11    "userdata_xval": X_val, "userdata_yval": y_val,
12    "userdata_xtest": X_test, "userdata_ytest": y_test,
13    "runcontrol_generations": 50,
14    "runcontrol_pop_size": 64,
15    "runcontrol_useparallel": False,
16    "functions_name": [["times", "minus", "plus", "divide"]],
17    "gene_max_genes": [5],
18    "tree_max_depth": [10],
19    "selection_elite_fraction": [0.05],
20    "selection_elite_fraction_ensemble": [0.05],
21    "selection_p_ensemble": [0.5],
22    "runcontrol_plotfitness": True,
23    "runcontrol_plotrankbest": True,
24    "runcontrol_plotsavefig": True,
25 }
26
27 gp = gpclassifier.initialize(**params)
28 gp.evolve()
29
30 # Write an HTML ensemble report for a chosen row (e.g., best on validation)
31 gp.report(ensemble_row=0) # creates an HTML file in the run's folder

```

## 5.6 Outputs: Plots and HTML Reports

**Plots.** If `runcontrol_plotfitness` or `runcontrol_plotrankbest` are enabled, the trainer writes PNGs (or configured formats) to the run's output directory. Typical figures include:

- Fitness vs. Generation (per population and/or aggregate),
- RankBest vs. Generation,
- (optional) RankAll vs. Generation.

**HTML report.** `gp.report(ensemble_row=...)` writes an HTML summary for the chosen ensemble row, typically including:

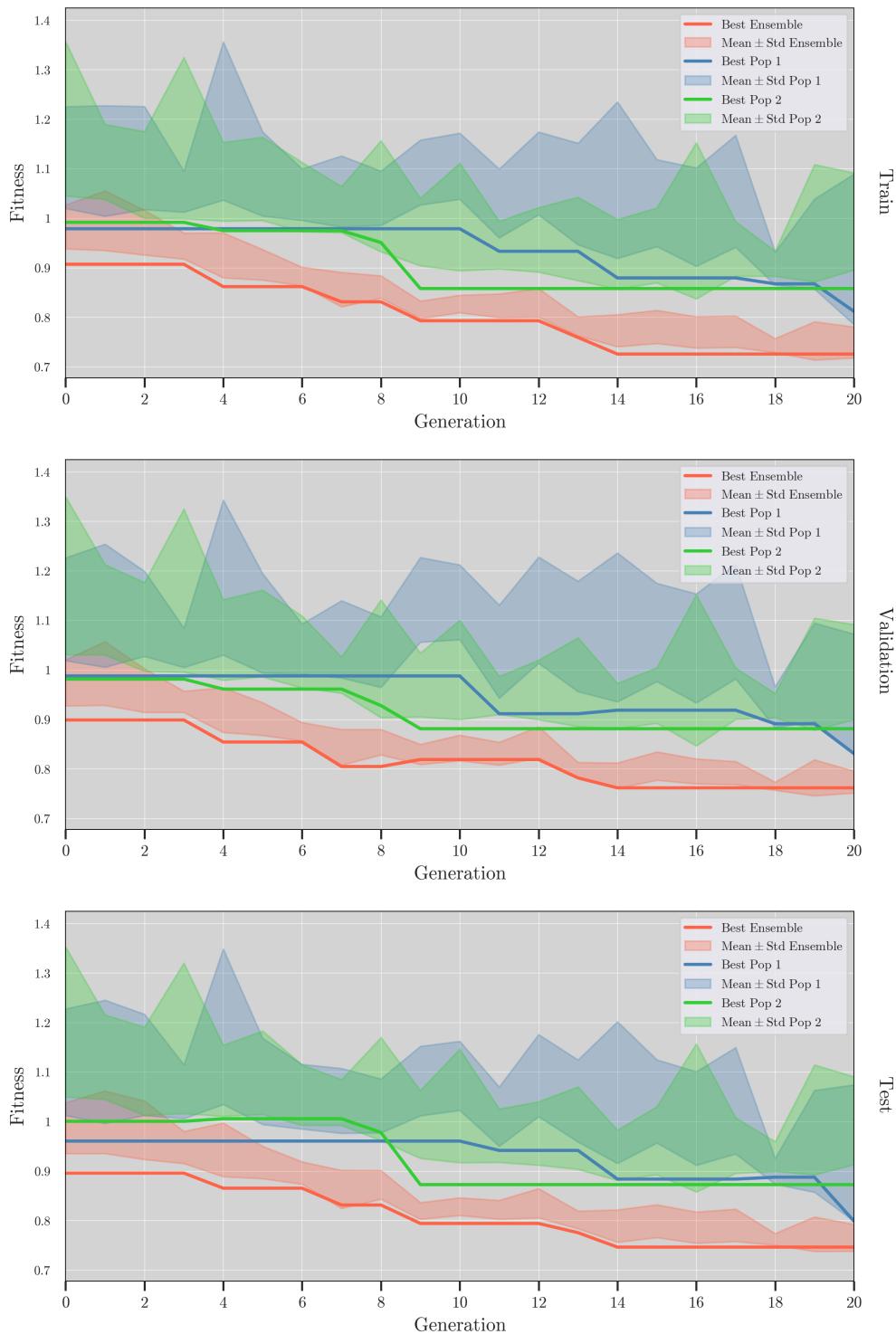
- per-population best trees (rendered),
- ensemble mixing weights and bias,
- validation/test metrics,
- run configuration snapshot.

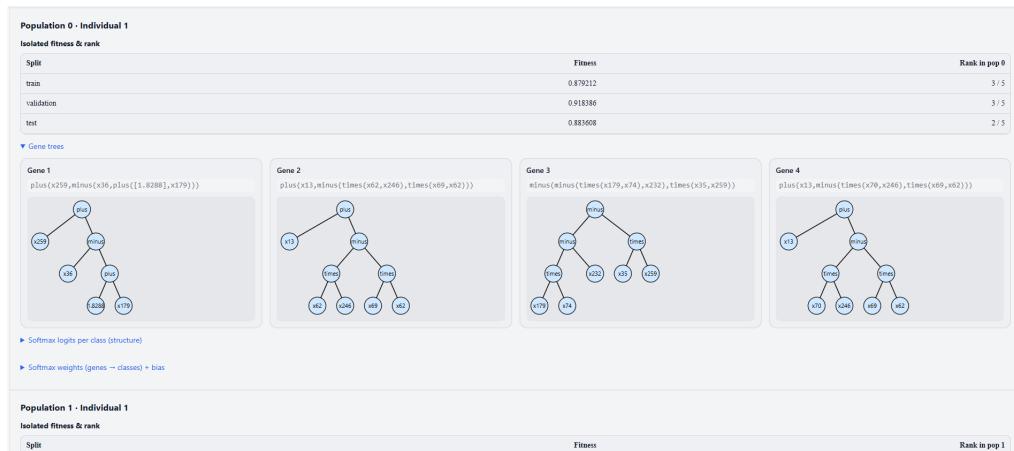
**Remark 5.2** (Assets to include). Please export and place the following under `Figure/` (or your chosen figures folder), then recompile:

1. `fig_clf_fitness.png`: Fitness vs. Generation from an actual run (referenced as Figure 5.1).
2. `fig_clf_report.png`: A cropped screenshot of the HTML report (referenced as Figure 5.2).

## 5.7 Tips & Common Pitfalls

- **Mismatched partitions.** When using `SPFPPartitioner`, ensure `userdata_pop_idx` has the same length as the number of populations implied by your configuration.
- **Function set arity.** The available primitives ("times", "minus", "plus", "divide", etc.) must match the arity expected by the tree builder; mixing incompatible sets can stall evolution.
- **Over-constrained structures.** Very small `tree_max_depth` or `gene_max_genes` can severely limit search; if fitness does not improve, relax these bounds.
- **Parallel evaluation.** If enabling `runcontrol_useparallel`, set `runcontrol_n_jobs` explicitly for predictable CPU usage.

**Figure 5.1:** Example *gpclassifier* Fitness vs. Generation.

**Figure 5.2:** Rendered gpclassifier HTML report.

# 6 gpregressor: Legacy Configuration & API

This chapter documents the `gpregressor` module: data requirements, training workflow, configuration keys (with representative defaults), how to integrate `SPFPPartitioner` partitions, and the produced artifacts (plots and HTML reports).

## 6.1 Overview & Data Requirements for `gpregressor`

**Task.** `gpregressor` implements multi-population Genetic Programming (GP) for supervised regression task. Multiple co-evolving populations produce scalar outputs that are linearly mixed by an ensemble head to form the final prediction. *Note that typed configurations are preferred over legacy dictionary for backward compatibility*

**Inputs.** Provide NumPy arrays:

- `userdata_xtrain`: shape  $(N_{\text{tr}}, D)$ ,
- `userdata_ytrain`: shape  $(N_{\text{tr}}, )$  with real targets,
- (optional) `userdata_xval`, `userdata_yval`,
- (optional) `userdata_xtest`, `userdata_ytest`.

If you use feature partitions from `SPFPPartitioner`, pass them via `userdata_pop_idx` (see Section 6.4).

## 6.2 Training Pipeline

A standard training call:

1. initializes populations (optionally restricted by `userdata_pop_idx`),

2. evolves for `runcontrol_generations` with selection, crossover, and mutation under structural limits,
3. fits an ensemble mixing head on population outputs (validation-aware when available),
4. evaluates and emits plots/reports/run-state.

## 6.3 Configuration Keys

### 6.3.1 User data & run metadata

**Table 6.1:** User data and metadata keys for gpregressor.

Key	Default	Description
<code>userdata_name</code>	"run"	Logical run name.
<code>userdata_xtrain</code>	—	Training features, ( $N_{tr}, D$ ).
<code>userdata_ytrain</code>	—	Training targets, ( $N_{tr},$ ).
<code>userdata_xval</code>	<code>None</code>	Validation features (optional).
<code>userdata_yval</code>	<code>None</code>	Validation targets (optional).
<code>userdata_xtest</code>	<code>None</code>	Test features (optional).
<code>userdata_ytest</code>	<code>None</code>	Test targets (optional).
<code>userdata_pop_idx</code>	<code>None</code>	Feature partitions (list of lists) from <code>SPFPPartitioner</code> .

### 6.3.2 Run control & parallelism

**Table 6.2:** Run control and parallel settings.

Key	Default	Description
<code>runcontrol_generations</code>	50	Number of generations.
<code>runcontrol_pop_size</code>	64	Individuals per population.
<code>runcontrol_useparallel</code>	<code>False</code>	Enable parallel fitness.
<code>runcontrol_n_jobs</code>	<code>None</code>	Workers (if parallel).
<code>runcontrol_batch_size</code>	<code>None</code>	Work chunk per worker.

### 6.3.3 Function sets, tree/gene constraints

**Table 6.3:** Structural configuration (function sets, depth, genes).

Key	Default	Description
functions_name	[["times", "minus", "plus", "divide"]]	Function set per population.
tree_max_depth	[10]	Max depth per population.
gene_max_genes	[5]	Max genes per chromosome (per population).

### 6.3.4 Selection & ensemble

**Table 6.4:** Selection and ensemble hyperparameters.

Key	Default	Description
selection_elite_fraction	[0.05]	Elitism per population.
selection_elite_fraction_ensemble	[0.05]	Elitism for ensemble stage.
selection_p_ensemble	[0.5]	Probability of ensemble candidate selection.

### 6.3.5 Plotting & reporting

**Table 6.5:** Plot and report toggles.

Key	Default	Description
runcontrol_plotfitness	True	Plot Fitness vs. Generation.
runcontrol_plotrankbest	True	Plot RankBest curve.
runcontrol_plotrankall	False	Plot RankAll curves.
runcontrol_plotsavefig	True	Save plots to disk.

**Remark 6.1** (About defaults). Defaults are representative and may vary by release. Consult the repository `README.md` and the source `gp_regress.py` for authoritative values.

## 6.4 Using SPFP Partitions

When using `SPFPPartitioner` (see Chapter 4), pass the resulting feature index groups via `userdata_pop_idx`. Ensure the number of groups equals the number of populations implied by your configuration:

**Listing 6.1:** Wiring SPFP partitions to gpregressor.

```

1 from genforge.gpregressor import gpregressor
2 from genforge.spfp.spfp_partition import SPFPPartitioner
3
4 spfp = SPFPPartitioner(n_groups=3, objective="micor", random_state=19)
5 spfp.fit(X, y)
6 findex = spfp.get_partitions()
7
8 params = {
9     "userdata_name": "reg_with_spfp",
10    "userdata_xtrain": X_train, "userdata_ytrain": y_train,
11    "userdata_xval": X_val, "userdata_yval": y_val,
12    "userdata_xtest": X_test, "userdata_ytest": y_test,
13    "userdata_pop_idx": findex, # <- map groups to populations
14    "runcontrol_generations": 50,
15    "runcontrol_pop_size": 64,
16    "functions_name": [["times", "minus", "plus", "divide"]],
17    "gene_max_genes": [5], "tree_max_depth": [10],
18    "selection_elite_fraction": [0.05],
19    "selection_elite_fraction_ensemble": [0.05],
20    "selection_p_ensemble": [0.5],
21    "runcontrol_plotfitness": True,
22    "runcontrol_plotrankbest": True,
23    "runcontrol_plotsavefig": True,
24 }
25 gp = gpregressor.initialize(**params)
26 gp.evolve()
27 gp.report(ensemble_row=0)

```

## 6.5 Training & Evaluation Workflow

Listing 6.2 shows a minimal end-to-end example. Adjust the structural constraints and function set to match your problem scale and variable types.

**Listing 6.2:** Minimal gpregressor workflow.

```

1 from genforge.gpregressor import gpregressor
2 import numpy as np
3
4 X_train, y_train = np.load("Xtr.npy"), np.load("ytr.npy")
5 X_val, y_val = np.load("Xva.npy"), np.load("yva.npy")
6 X_test, y_test = np.load("Xte.npy"), np.load("yte.npy")
7
8 params = {
9     "userdata_name": "demo_reg",
10    "userdata_xtrain": X_train, "userdata_ytrain": y_train,
11    "userdata_xval": X_val, "userdata_yval": y_val,
12    "userdata_xtest": X_test, "userdata_ytest": y_test,
13    "runcontrol_generations": 50,
14    "runcontrol_pop_size": 64,
15    "runcontrol_useparallel": False,
16    "functions_name": [["times", "minus", "plus", "divide"]],
17    "gene_max_genes": [5],
18    "tree_max_depth": [10],
19    "selection_elite_fraction": [0.05],
20    "selection_elite_fraction_ensemble": [0.05],
21    "selection_p_ensemble": [0.5],
22    "runcontrol_plotfitness": True,
23    "runcontrol_plotrankbest": True,
24    "runcontrol_plotsavefig": True,
25 }
26
27 gp = gpregressor.initialize(**params)
28 gp.evolve()
29
30 # Write an HTML ensemble report for a chosen row (e.g., best on validation)
31 gp.report(ensemble_row=0)
```

## 6.6 Outputs: Plots and HTML Reports

**Plots.** If `runcontrol_plotfitness` and/or `runcontrol_plotrankbest` are enabled, the trainer saves figures to the run's output directory. Common plots:

- Fitness vs. Generation,
- RankBest vs. Generation,
- (optional) RankAll vs. Generation.

**HTML report.** `gp.report(ensemble_row=...)` writes a self-contained HTML summary with:

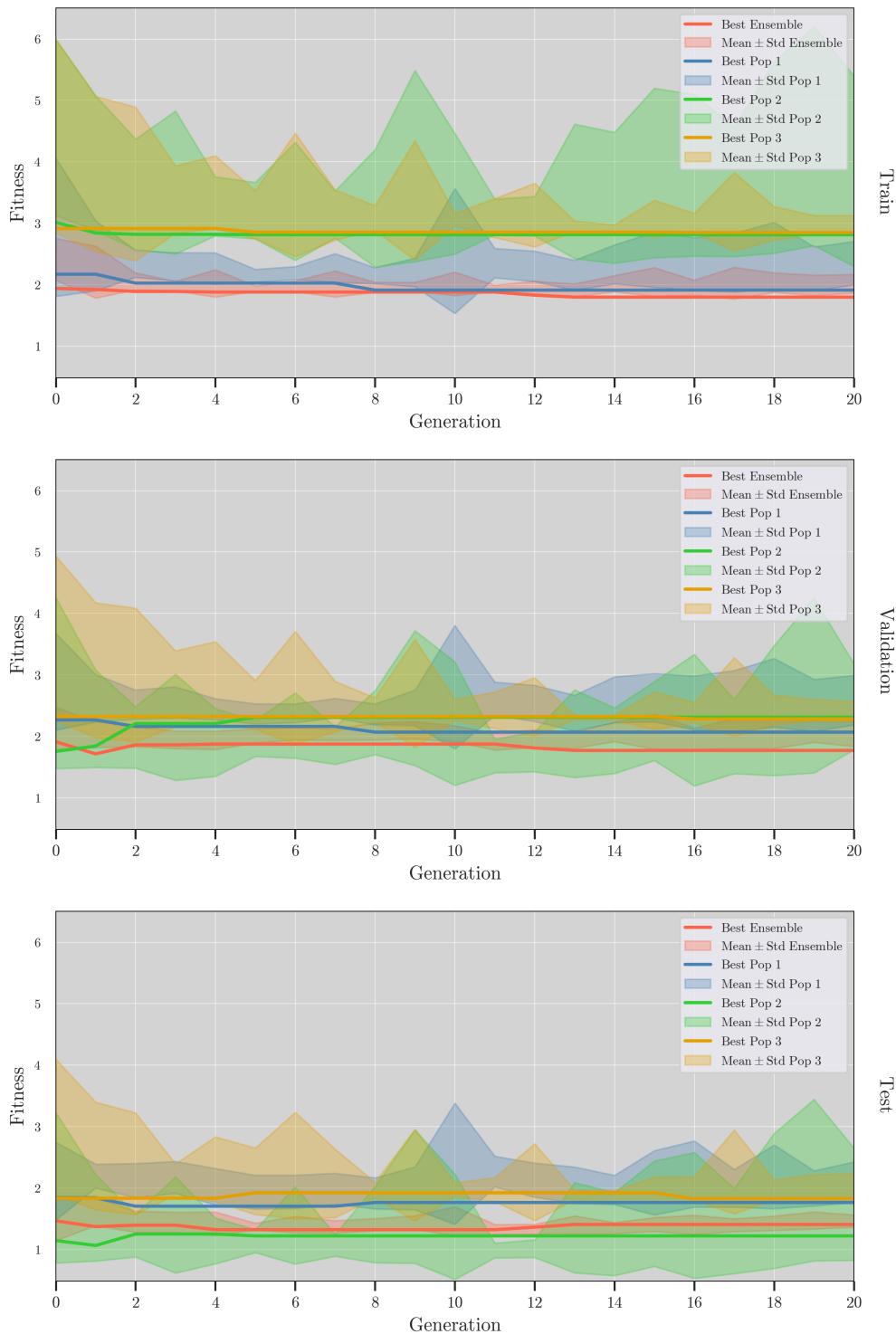
- best trees per population (rendered),
- ensemble mixing weights and bias,
- validation/test metrics (e.g., MAE, RMSE as available),
- configuration snapshot.

**Remark 6.2** (Assets to include). Please export and place the following under `Figure/`, then recompile:

1. `fig_reg_fitness.png`: Fitness vs. Generation from a real run (for Figure 6.1).
2. `fig_reg_report.png`: Cropped screenshot of the regression HTML report (for Figure 6.2).

## 6.7 Tips & Common Pitfalls

- **Scale of targets.** Extremely large or small target magnitudes can affect numerical stability. Consider rescaling `userdata_ytrain` (and corresponding splits) if necessary.
- **Partitions vs. populations.** When using `SPFPPartitioner`, ensure the number of groups equals the number of populations; otherwise, features may be unassigned or over-constrained.
- **Search budget.** If fitness stagnates, increase `runcontrol_generations`, `runcontrol_pop_size`, or relax `tree_max_depth/gene_max_genes`.
- **Parallel evaluation.** When enabling `runcontrol_useparallel`, set `runcontrol_n_jobs` explicitly for predictable CPU usage.

**Figure 6.1:** Example *gpregressor* Fitness vs. Generation.

<b>SFRC_Example: chosen ensemble #0</b>				
This report shows the chosen ensemble, its member individuals per population, and each member's gene trees. Node labels are the exact function names from the GenForge run (e.g., plus, minus, times, divide). For regression, we display the learned linear head weights (genes → target) with an explicit <i>bias</i> column when available, ensemble mixing weights, and fitness values (with ranks).				
<b>Ensemble mixing weights</b>				
Target	w[pop=0]	w[pop=1]	w[pop=2]	bias
0	0.803166	0.201508	0.0520948	-0.147363

<b>Ensemble fitness &amp; rank</b>		
Split	Fitness	Rank among ensembles
train	1.83549	2 / 10
validation	2.03791	5 / 10
test	1.44816	6 / 10

<b>Population 0 - Individual 0</b>		
<b>Isolated fitness &amp; rank</b>		
Split	Fitness	Rank in pop 0
train	1.90968	3 / 10
validation	2.06678	1 / 10

**Figure 6.2:** Rendered gpregressor HTML report.

# 7 SPFPPartitioner: Configuration & API

This chapter documents the `SPFPPartitioner` module as shipped in `GenForge`: its public constructor, methods, expected inputs/outputs, and how to integrate the produced feature partitions (`findex`) with `gpclassifier` and `gpregressor`.

## 7.1 Overview

**Goal.** `SPFPPartitioner` (Semantic Preserving Feature Partitioning) groups the  $D$  input features into  $P$  disjoint (or near-disjoint) index sets, each intended to feed one GP population. The grouping is driven by entropy-based relevance to the label and by pairwise interaction/redundancy controls, following the exact mathematics in Chapter 4.

**Inputs.** Provide `numpy.ndarray` arrays:

- `x`: shape  $(N, D)$ ; features are internally standardised (zero mean, unit variance) with the built-in `StandardScaler`.
- `y`: shape  $(N, )$ ; labels are internally encoded by the built-in `LabelEncoder` (classification) or used as-is (regression targets).

**Outputs.** After fitting, `SPFPPartitioner` exposes `partitions_` (list of  $P$  lists of feature indices) and optional diagnostics (per-feature MI/correlation tables, group entropies, timings).

## 7.2 Constructor Parameters

**Remark 7.1.** All MI-related terms are normalised by  $H(y)$ , and correlation uses the absolute Pearson coefficient on the *continuous, standardised* `x` against the encoded `y`, exactly as in the reference script.

**Table 7.1:** Public constructor of `SPFPPartitioner`. All defaults mirror the reference implementation.

Parameter	Default	Meaning
<code>n_groups</code>	5	Number of groups $P$ to construct.
<code>n_bins</code>	100	Equal-distance discretisation bins for entropy ops.
<code>objective</code>	"micor"	Primary score in gated branch: "mi", "cor", or "micor".
<code>criterion_add</code>	0.999	$\alpha$ in forward stopping rule (Section 4.3).
<code>criterion_rem</code>	0.999	$\beta$ in optional backward pruning (Section 4.4).
<code>growth_magnitude</code>	1.01	$\gamma$ cap during growth-gated search.
<code>perfs</code>	0.10	$\pi$ : minimum ratio of $D$ per group.
<code>remp</code>	0.60	$\rho$ : random withholding rate between groups.
<code>override</code>	False	If <code>True</code> , pick by composite objective directly.
<code>backward</code>	False	Enable pruning after forward selection.
<code>floor_decimals</code>	4	Decimal precision for threshold floors.
<code>random_state</code>	None	Seed for reproducibility.
<code>verbose</code>	1	Verbosity level (e.g., progress prints).

## 7.3 Public Methods

Some implementations also expose `get_params()` / `set_params(**kw)` for convenience; consult the class docstring in `spfp_partition.py`.

## 7.4 Minimal Usage

**Listing 7.1:** Compute feature partitions with `SPFPPartitioner`.

```

1 from genforge.spfp.spfp_partition import SPFPPartitioner
2 import numpy as np
3
4 # Load arrays
5 X = np.load("X.npy") # shape (N, D)
6 y = np.load("y.npy") # shape (N,)
7
8 spfp = SPFPPartitioner(
9     n_groups=5,

```

**Table 7.2:** Core methods of SPFPPartitioner.

Signature	Description
<code>fit(X, y)</code>	Runs standardisation/encoding, per-feature scoring, group-wise forward selection (and optional backward pruning), then inter-group withholding. Populates <code>partitions_</code> and diagnostics.
<code>partition()</code>	Builds groups using the statistics prepared by <code>fit</code> . Typically called internally by <code>fit</code> ; exposed for advanced workflows.
<code>fit_transform(X, y)</code>	Convenience wrapper: <code>fit</code> then return a lightweight result (e.g., <code>partitions_</code> ).
<code>get_partitions()</code>	Returns the list of groups: <code>[[idx...], [idx...], ...]</code> .
<code>transform(X)</code>	Optional helper to select columns per group; returns a list of views or a dict <code>{group_id: X[:, idx]}</code> .
<code>get_diagnostics()</code>	Returns a dictionary of MI/correlation tables and group entropy summaries, useful for reporting.
<code>get_binning_edges()</code>	Returns the per-feature discretisation edges used by the joint histogram.

```

10 n_bins=100,
11 objective="micor", # "mi" / "cor" / "micor"
12 criterion_add=0.999,
13 criterion_rem=0.999,
14 growth_magnitude=1.01,
15 perfs=0.10,
16 remp=0.60,
17 override=False,
18 backward=False,
19 random_state=13,
20 verbose=1
21 )
22
23 spfp.fit(X, y)
24 findex = spfp.get_partitions() # list of lists with feature indices
25
26 for g, idx in enumerate(findex):
27     print(f"Group {g}: {idx}")

```

**Diagnostics (optional).****Listing 7.2:** Inspect SPFP diagnostics.

```

1 diag = spfp.get_diagnostics()
2 # e.g., per-feature MI/correlation table, group entropies, timings:
3 mi_table = diag.get("feature_table") # implementation-dependent
4 group_info = diag.get("group_summaries")
5 print("Keys:", list(diag.keys()))

```

## 7.5 Wiring Partitions to GP

Use the produced `findex` to constrain populations in your GP run:

**Listing 7.3:** Using SPFP partitions with `gpclassifier` (analogous for `gpregressor`).

```

1 from genforge(gpclassifier import gpclassifier
2 from genforge.spfp.spfp_partition import SPFPPartitioner
3 import numpy as np
4
5 X_train = np.load("X_train.npy")
6 y_train = np.load("y_train.npy")
7 X_val = np.load("X_val.npy")
8 y_val = np.load("y_val.npy")
9 X_test = np.load("X_test.npy")
10 y_test = np.load("y_test.npy")
11
12 # Build 3 groups for 3 populations
13 spfp = SPFPPartitioner(n_groups=3, objective="micor", random_state=7)
14 spfp.fit(X_train, y_train)
15 findex = spfp.get_partitions()
16
17 params = {
18     "userdata_name": "clf_with_spfp",
19     "userdata_xtrain": X_train, "userdata_ytrain": y_train,
20     "userdata_xval": X_val, "userdata_yval": y_val,
21     "userdata_xtest": X_test, "userdata_ytest": y_test,
22     "userdata_pop_idx": findex, # <- map groups to populations
23     "runcontrol_generations": 50,
24     "runcontrol_pop_size": 64,

```

```

25 "functions_name": [["times", "minus", "plus", "divide"]],
26 "gene_max_genes": [5], "tree_max_depth": [10],
27 "selection_elite_fraction": [0.05],
28 "selection_elite_fraction_ensemble": [0.05],
29 "selection_p_ensemble": [0.5],
30 "runcontrol_plotfitness": True,
31 "runcontrol_plotrankbest": True,
32 "runcontrol_plotsavefig": True,
33 }
34
35 gp = gpclassifier.initialize(**params)
36 gp.evolve()
37 gp.report(ensemble_row=0)

```

**Remark 7.2.** Ensure the number of groups equals the number of populations implied by your GP configuration; otherwise, some populations may receive empty feature sets.

## 7.6 Notes & Gotchas

- **Scaling and encoding.** `SPFPPartitioner` applies its own feature standardisation and label encoding; pass raw `X`, `y` to it.
- **Bin edges.** Training-set binning edges are re-used consistently for any test-time diagnostics you request, to keep entropies comparable.
- **Objective choice.** `objective` only affects a specific selection branch (when coverage thresholds are met but the minimum size is not). Otherwise, growth-gated search or override governs selection.
- **Reproducibility.** Set `random_state` to stabilise stream ordering and the inter-group withholding step.



# 8 Plots & HTML Reports

GenForge ships with lightweight plotting utilities and an HTML report mechanism to inspect the final ensemble and its member programs. This chapter explains how to enable/save plots, what each figure means, and how to generate the HTML report from `gpclassifier` and `gpregressor`.

## 8.1 Enabling and Saving Plots

Both `gpclassifier` and `gpregressor` support the following run-control switches (names match the code in `gp_class.py` and `gp_regress.py`). If a key is not present for your estimator, it is silently ignored.

**Table 8.1:** Common plot flags and output controls.

Parameter	Meaning	Typical values
<code>runcontrol_plotfitness</code>	Enable Fitness vs. Generation plot	True/False
<code>runcontrol_plotrankbest</code>	Enable RankBest vs. Generation plot	True/False
<code>runcontrol_plotrankall</code>	Enable RankAll vs. Generation plot	True/False
<code>runcontrol_plotsavefig</code>	Save figures to disk	True/False
<code>runcontrol_plotformat</code>	Image formats to save	e.g., ["png", "pdf"]
<code>runcontrol_plotfolder</code>	Output folder for plots	e.g., "./plots/"
<code>runcontrol_plotbackend</code>	Matplotlib backend (regressor supports)	"Agg", "TkAgg"
<code>runcontrol_plotlive</code>	Live drawing during run (regressor)	True/False

Example configuration snippet (classification) activating all three plots and saving to `./plots/` in `png` and `pdf`:

**Listing 8.1:** Plot configuration example for `gpclassifier`.

```
1 params.update({
2     "runcontrol_plotfitness": True,
```

```

3   "runcontrol_plotrankbest": True,
4   "runcontrol_plotrankall": True,
5   "runcontrol_plotsavefig": True,
6   "runcontrol_plotformat": ["png", "pdf"],
7   "runcontrol_plotfolder": "./plots/"
8 })

```

The same flags apply to `gpregressor`; additionally, the regressor honors `runcontrol_plotbackend` and `runcontrol_plotlive` if present.

## 8.2 Fitness vs. Generation

The Fitness curve tracks the best (or aggregate) fitness value per generation. For minimisation runs (the default, `runcontrol_minimisation = True`), downward is better; for maximisation, upward improves.

**Tips.** Sudden jumps typically indicate successful crossover/mutation; stagnation may suggest increasing `runcontrol_pop_size`, relaxing constraints, or injecting diversity.

## 8.3 RankBest vs. Generation

RankBest reports the rank (1 is best) of the best individual across generations within each population or globally (implementation-dependent). It is a compact progress signal that is less sensitive to fitness scale.

## 8.4 RankAll vs. Generation

RankAll visualises the distribution of ranks over time (e.g., median and spread). It complements RankBest by revealing whether the whole population improves or only a few elites do.

## 8.5 HTML Ensemble Report

After evolution, call `report` to emit an HTML page that summarises the chosen ensemble row (e.g., best on the validation split), its mixing weights, and member programs.

**Listing 8.2:** Emit an HTML report from a trained run.

```
1 # Classification
```

```

2 gp = gpclassifier.initialize(**params)
3 gp.evolve()
4 gp.report(ensemble_row=0) # writes an HTML file for row 0
5
6 # Regression
7 gr = gpregressor.initialize(**params)
8 gr.evolve()
9 gr.report(ensemble_row=0) # same idea

```

**What is included.** The exact layout depends on your version, but typically contains:

- **Ensemble mixing weights** (including bias term), mapping population outputs to the final prediction.
- **Member programs** (trees) with human-readable expressions.
- **Per-split metrics** (training/validation/test).
- **Run metadata** (seeds, wall time, configuration summary).

**Location and naming.** By default, reports are written alongside run artifacts (plots/logs). If your build has a configurable output root, set it via the run-control keys used in your environment (e.g., `runcontrol_plotfolder` for figures).

## 8.6 Interpreting the Ensemble

**Mixing weights.** Large-magnitude weights highlight populations contributing most to the final prediction. If you use `SFPPPartitioner`, you can attribute contributions back to feature groups by mapping each population to its `index`.

**Program inspection.** Member trees often reveal structure such as repeated sub-expressions, protected divisions, or domain-specific interactions (e.g., `times`, `minus`). Use these clues to refine the function set or add domain constraints.

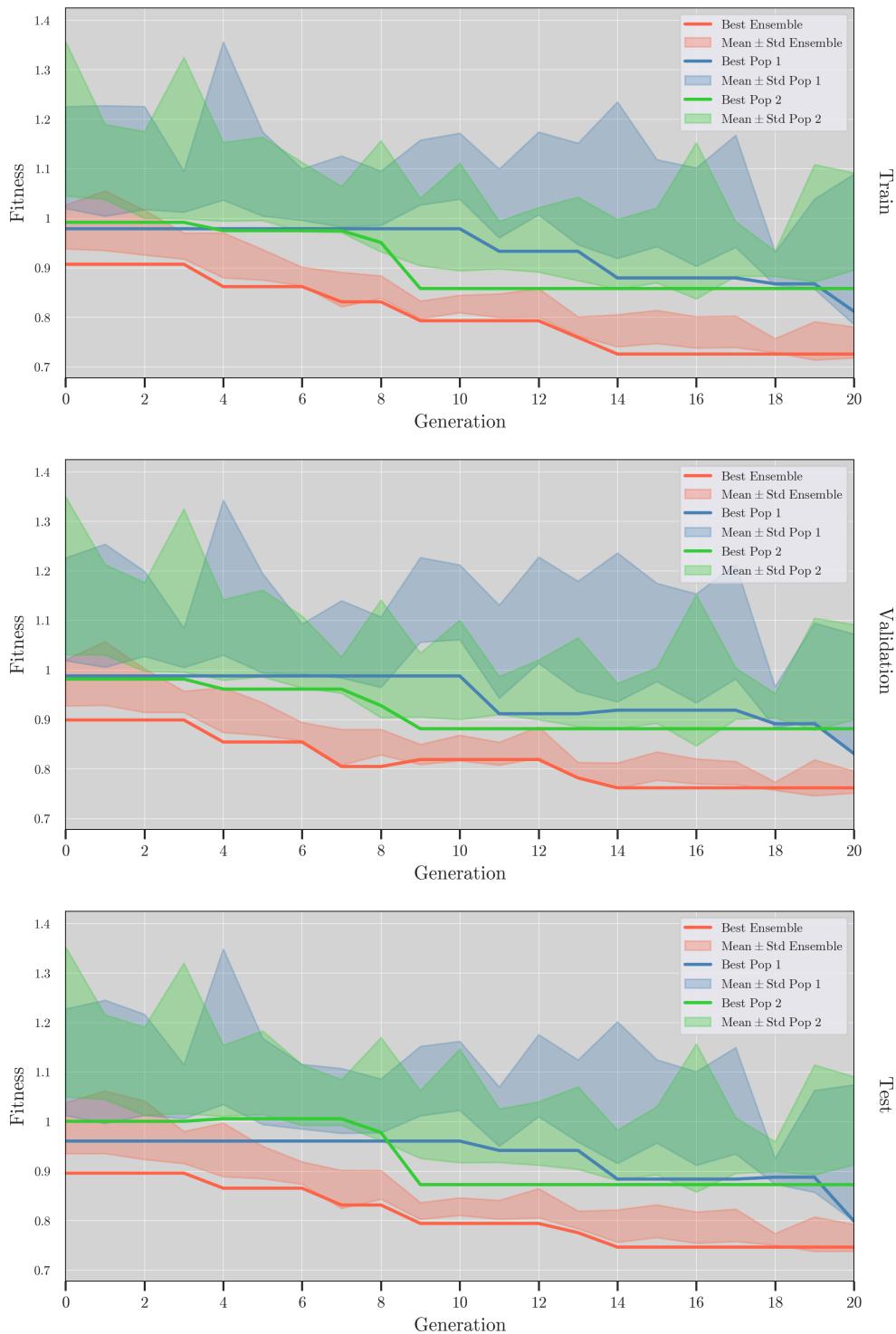
## 8.7 Assets to Include for This Chapter

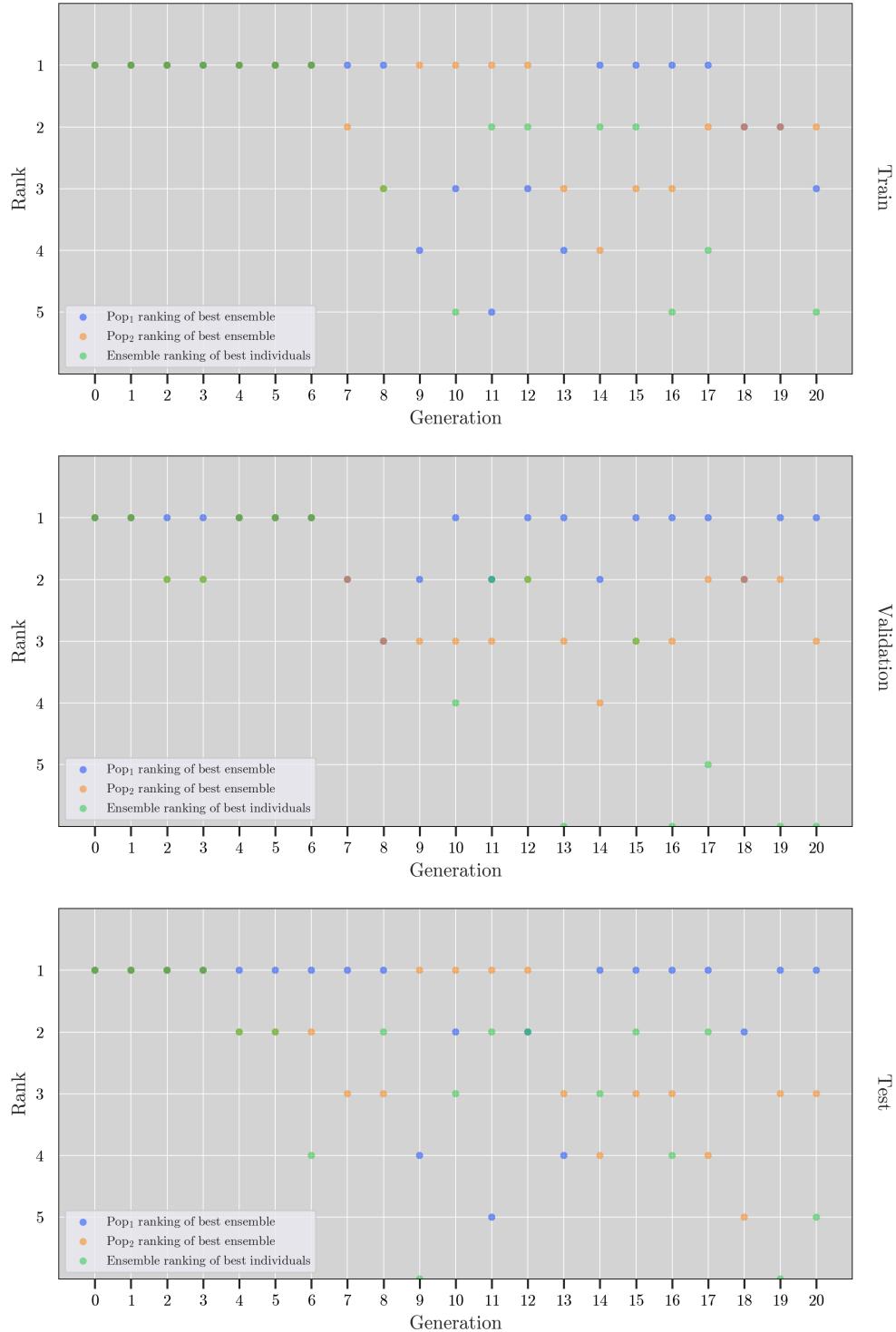
**Remark 8.1** (Asset checklist). Please export and place the following files under `Figure/` (or update paths accordingly); then recompile to populate Figures 8.1 to 8.4.

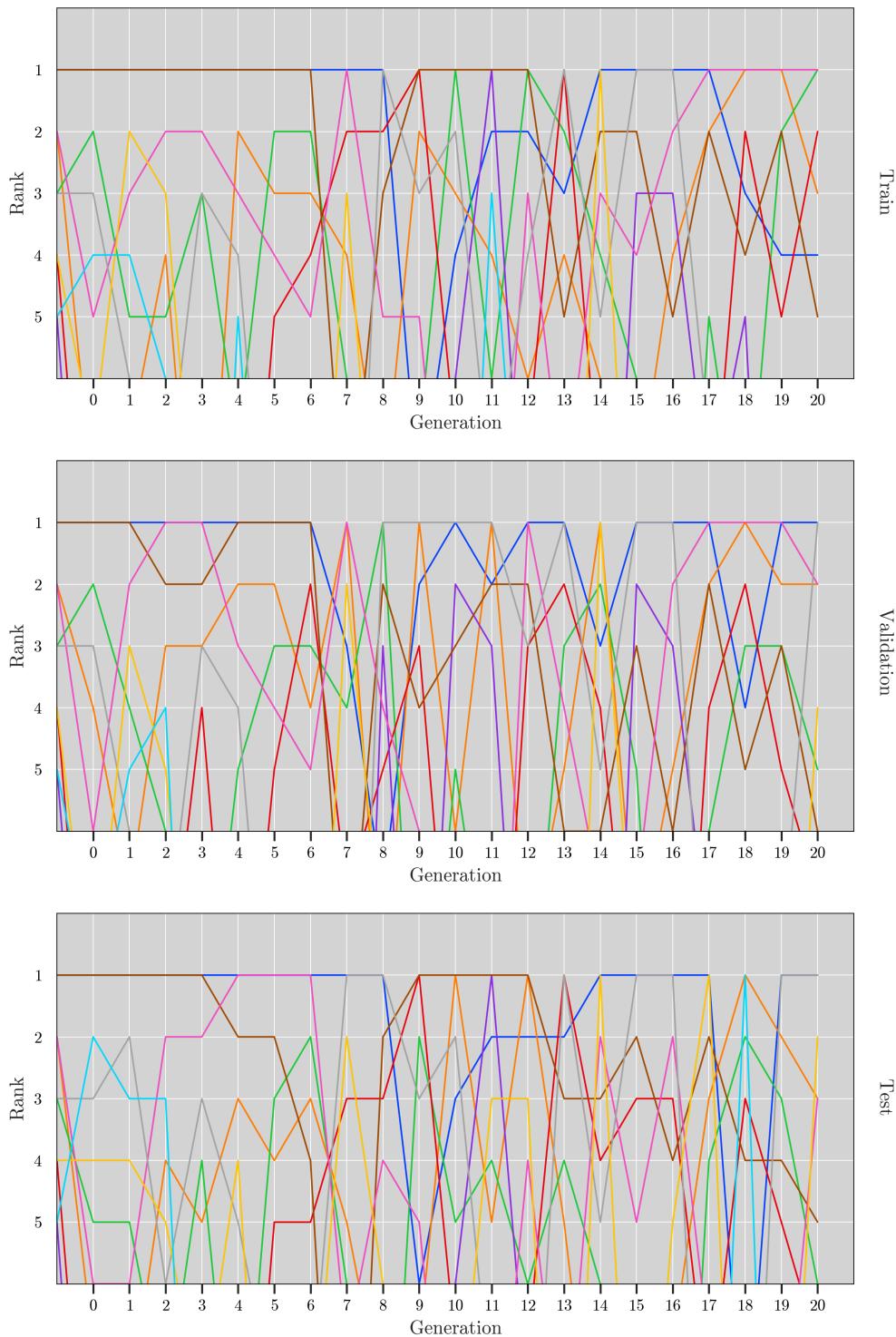
1. `fig_fitness_vs_generation.png`: Fitness plot (Figure 8.1).
2. `fig_rankbest_vs_generation.png`: RankBest plot (Figure 8.2).
3. `fig_rankall_vs_generation.png`: RankAll plot (Figure 8.3).
4. `fig_html_report_screenshot.png`: Screenshot of a generated HTML report (Figure 8.4)—capture the ensemble weights and a couple of member trees.

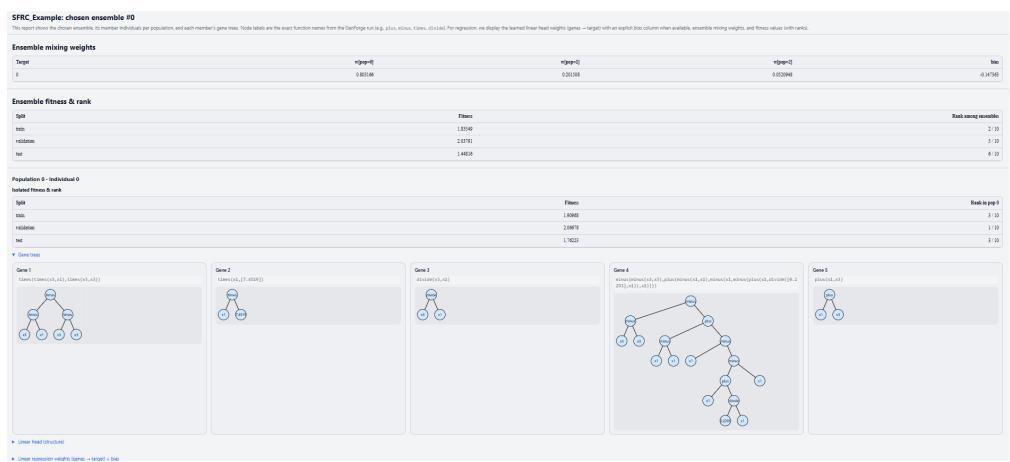
## 8.8 Troubleshooting

- **Figures do not appear.** Ensure `runcontrol_plotsavefig = True`, the folder in `runcontrol_plotfolder` exists, and your process has write permission.
- **Blank plots on a headless server.** Set `runcontrol_plotbackend = "Agg"` (supported by the regressor) and disable live drawing.
- **HTML report not created.** Confirm you called `report` after training and that the `ensemble_row` you requested exists.
- **Stale ToC/LoF/LoT entries.** Compile twice. If labels still show `???`, verify file paths in `\graphicspath` and figure filenames.

**Figure 8.1:** Fitness vs. Generation for a representative run.

**Figure 8.2:** RankBest vs. Generation (lower is better).

**Figure 8.3:** RankAll vs. Generation showing median and dispersion.

**Figure 8.4:** HTML report (example) summarising an ensemble row.

# 9 Reproducibility & Experiment Management

This chapter outlines practical steps to make GenForge runs repeatable and easy to manage at scale: seeding, controlling parallelism and BLAS threads, keeping data splits consistent, organising artifacts, and comparing runs.

## 9.1 Goals and Scope

Reproducibility in evolutionary methods is nuanced: identical seeds and configurations should yield comparable outcomes, but parallelism, floating-point reductions, and BLAS libraries may introduce small deviations. The guidance below focuses on *practical determinism*—repeatable enough for fair comparisons and ablation studies.

## 9.2 Seeding and Determinism

**Python/NumPy.** Always seed the Python and NumPy RNGs at process start:

**Listing 9.1:** Seeding Python and NumPy.

```
1 import os, random, numpy as np
2 os.environ["PYTHONHASHSEED"] = "0" # make hashed iterations stable
3 random.seed(1337)
4 np.random.seed(1337)
```

**SPFPPartitioner.** Set `random_state` in the `SPFPPartitioner` constructor for a stable partition:

**Listing 9.2:** Seeding SPFP.

```
1 from genforge.spfp.spfp_partition import SPFPPartitioner
2 spfp = SPFPPartitioner(random_state=1337, n_groups=5, n_bins=100, objective="micor")
```

**gpclassifier / gpregressor.** If your build exposes a run seed (e.g., `userdata_seed` or similar in your environment), set it in `params`. When no explicit run seed exists, use the global RNG seeds (Listing 9.1) and keep `runcontrol_useparallel` disabled for strict determinism:

**Listing 9.3:** Minimal deterministic GP run.

```

1 from genforge.gpclassifier import gpclassifier
2 params = {
3     "userdata_name": "deterministic_demo",
4     # "userdata_seed": 1337, # if your build supports one
5     "runcontrol_useparallel": False,
6     "runcontrol_generations": 10,
7     "runcontrol_pop_size": 32,
8     "functions_name": [["times", "minus", "plus", "divide"]],
9     "gene_max_genes": [5], "tree_max_depth": [10],
10    "selection_elite_fraction": [0.05],
11    "selection_elite_fraction_ensemble": [0.05],
12    "selection_p_ensemble": [0.5],
13 }
14 gp = gpclassifier.initialize(**params)
15 gp.evolve()
```

## 9.3 Controlling BLAS & Threading

Numeric libraries (OpenBLAS, MKL, BLIS, numexpr) spawn their own threads, which can affect both speed and determinism. Set thread caps *before* importing NumPy/Scipy:

**Table 9.1:** Useful environment variables for threading control.

Variable	Recommended setting (example)
<code>OMP_NUM_THREADS</code>	1 (strict) or a small integer
<code>OPENBLAS_NUM_THREADS</code>	1
<code>MKL_NUM_THREADS</code>	1
<code>VECLIB_MAXIMUM_THREADS</code>	1 (Accelerate, macOS)
<code>NUMEXPR_NUM_THREADS</code>	1

Example launcher (Linux/macOS):

**Listing 9.4:** Launch with controlled threads.

```
1 OMP_NUM_THREADS=1 OPENBLAS_NUM_THREADS=1 MKL_NUM_THREADS=1 \
2 VECLIB_MAXIMUM_THREADS=1 NUMEXPR_NUM_THREADS=1 \
3 python run_experiment.py
```

If you rely on parallel fitness evaluation, prefer a *single BLAS thread* with multiple Python workers via `runcontrol_useparallel` and (if available) `runcontrol_n_jobs`. This reduces contention.

## 9.4 Data Splits and Consistency

Use the same training split for both `SPFPPartitioner` partitioning and GP evolution to avoid leakage.

For classification with imbalanced classes, stratify externally when preparing `X_train/y_train`. For time series, avoid shuffling across temporal boundaries.

**Remark 9.1** (One source of truth for splits). Persist `X_train.npy`, `y_train.npy`, `X_val.npy`, `y_val.npy`, `X_test.npy`, `y_test.npy` and reuse them across all runs. Document the split provenance (random seed, date, filter criteria).

## 9.5 Configuration Management

Keep run configurations as plain `dict` objects and persist them (JSON/YAML). This ensures you can regenerate the same run later.

**Listing 9.5:** Save and reload run configuration.

```
1 import json, pathlib
2 cfg_path = pathlib.Path("runs/demo_clf/config.json")
3 cfg_path.parent.mkdir(parents=True, exist_ok=True)
4 with cfg_path.open("w") as f:
5     json.dump(params, f, indent=2, sort_keys=True)
6
7 # ... later
8 with cfg_path.open() as f:
9     params2 = json.load(f)
```

**Table 9.2:** Suggested minimal keys to track per run (extend as needed).

Key	Purpose
userdata_name	Human-readable run name
runcontrol_generations	Evolution length
runcontrol_pop_size	Population size
runcontrol_useparallel	Enable/disable parallel evaluation
functions_name	Function set per population
gene_max_genes, tree_max_depth	Structural limits
selection_elite_fraction	Elitism fraction
selection_p_ensemble	Ensemble selection probability
userdata_pop_idx	(Optional) SPFPPartitioner <i>findex</i> mapping
random_state	(Where applicable) run seed

## 9.6 Run Naming and Artifact Layout

Adopt a predictable folder layout to store plots, HTML reports, and configs. A simple pattern is `runs/<name>/...:`

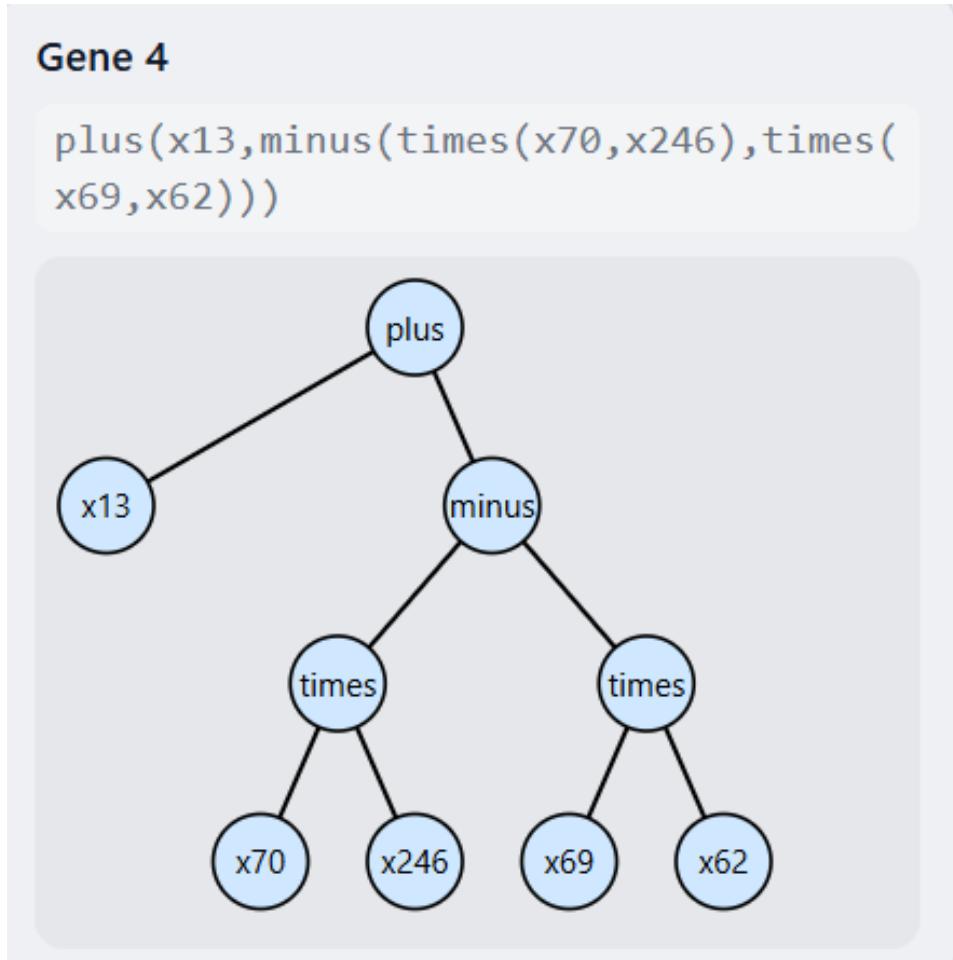
**Listing 9.6:** Example artifact tree (customise to taste).

```

1 runs/
2   demo_clf/
3     config.json
4   plots/
5     fitness.png
6     rankbest.png
7     rankall.png
8   reports/
9     ensemble_row0.html
10  logs/
11    stdout.txt
12    metrics.csv

```

If your build already writes plots/reports to a default location (e.g., `runcontrol_plotfolder`), point it to `runs/<name>/plots/` for consistency.



**Figure 9.1:** Example `runs/` directory with plots, reports, and logs.

## 9.7 Logging Metrics and Comparing Runs

Log core metrics to a CSV during/after training (validation accuracy, test MAE, wall time). A minimal schema:

**Listing 9.7:** Append a metric row to CSV.

```

1 import csv, time, pathlib
2 csvp = pathlib.Path("runs/demo_clf/metrics.csv")
3 csvp.parent.mkdir(parents=True, exist_ok=True)
4 with csvp.open("a", newline="") as f:

```

```

5   w = csv.writer(f)
6   if f.tell() == 0:
7       w.writerow(["run", "split", "metric", "value", "step", "timestamp"])
8   w.writerow(["demo_clf", "val", "accuracy", 0.842, 50, int(time.time())])

```

Aggregate multiple runs in a notebook or script to produce comparison tables and stability plots.

## 9.8 Strict vs. Practical Determinism

- **Strict (bit-wise) determinism:** Disable parallelism (`runcontrol_useparallel = False`), pin BLAS threads to 1 (Table 9.1), avoid nondeterministic reductions, and run on the same machine/OS with identical library versions.
- **Practical determinism:** Allow parallel evaluation but keep BLAS threads at 1 and fix seeds. Expect tiny variations in tie-breakers.

## 9.9 End-to-End Example

A compact script that (i) seeds everything, (ii) computes `SPFPPartitioner` partitions, (iii) runs a GP classifier with those partitions, and (iv) writes plots/reports.

**Listing 9.8:** Seeded pipeline: SPFP + GPClassifier.

```

1 import os, random, json, numpy as np, pathlib
2 from genforge.spfp.spfp_partition import SPFPPartitioner
3 from genforge.gpclassifier import gpclassifier
4
5 # --- seeds and threads
6 os.environ.update({
7     "PYTHONHASHSEED": "0",
8     "OMP_NUM_THREADS": "1", "OPENBLAS_NUM_THREADS": "1",
9     "MKL_NUM_THREADS": "1", "VECLIB_MAXIMUM_THREADS": "1",
10    "NUMEXPR_NUM_THREADS": "1"
11 })
12 random.seed(1337); np.random.seed(1337)
13
14 # --- data (replace with your arrays)
15 X_train = np.load("data/X_train.npy")

```

```

16 y_train = np.load("data/y_train.npy")
17 X_val = np.load("data/X_val.npy")
18 y_val = np.load("data/y_val.npy")
19 X_test = np.load("data/X_test.npy")
20 y_test = np.load("data/y_test.npy")
21
22 # --- SPFP partitions
23 spfp = SPFPPartitioner(n_groups=4, n_bins=100, objective="micor",
24                         perfs=0.10, remp=0.60, random_state=1337, verbose=1)
25 spfp.fit(X_train, y_train)
26 findex = spfp.get_partitions() # [[feature indices per group], ...]
27
28 # --- GP configuration
29 params = {
30     "userdata_name": "spfp_gp_demo",
31     "userdata_xtrain": X_train, "userdata_ytrain": y_train,
32     "userdata_xval": X_val, "userdata_yval": y_val,
33     "userdata_xtest": X_test, "userdata_ytest": y_test,
34     "userdata_pop_idx": findex, # wire features to populations
35     "runcontrol_generations": 30,
36     "runcontrol_pop_size": 48,
37     "runcontrol_useparallel": False, # stricter determinism
38     "runcontrol_plotfitness": True,
39     "runcontrol_plotrankbest": True,
40     "runcontrol_plotsavefig": True,
41     "runcontrol_plotfolder": "runs/spfp_gp_demo/plots/",
42     "runcontrol_plotformat": ["png"],
43     "functions_name": [["times", "minus", "plus", "divide"]],
44     "gene_max_genes": [5], "tree_max_depth": [10],
45     "selection_elite_fraction": [0.05],
46     "selection_elite_fraction_ensemble": [0.05],
47     "selection_p_ensemble": [0.5],
48 }
49 runroot = pathlib.Path("runs/spfp_gp_demo")
50 (runroot / "plots").mkdir(parents=True, exist_ok=True)
51 with (runroot / "config.json").open("w") as f:
52     json.dump(params, f, indent=2)
53

```

```
54 # --- train and report
55 gp = gpclassifier.initialize(**params)
56 gp.evolve()
57 gp.report(ensemble_row=0)
```

## 9.10 Troubleshooting

- **Different results across machines.** Pin versions in `requirements.txt` or `pyproject.toml`; align BLAS libraries and thread settings (Table 9.1).
- **Plots missing.** Ensure `runcontrol_plotsavefig = True`, and `runcontrol_plotfolder` exists and is writable.
- **HTML not generated.** Call `report(ensemble_row=...)` after training and verify the requested row index exists.
- **SPFP partitions change run-to-run.** Set `random_state` in `SPFPPartitioner` and reuse the same training split.

# 10 Advanced Usage & Extensibility

This chapter collects power-user patterns and extension points for `GenForge`: custom primitive function sets, feature-group wiring via `SPFPPartitioner`, constants configuration, selection/ensemble controls, parallel evaluation, and practical performance tips.

## 10.1 Wiring SPFPPartitioner to Multi-Population GP

Given  $P$  feature groups produced by `SPFPPartitioner` (see Chapter 4), map them to GP populations by passing the per-group index lists through `userdata_pop_idx`. The number of populations equals the number of groups.

**Listing 10.1:** Using SPFPPartitioner partitions with a multi-population run.

```
1 from genforge.spfp.spfp_partition import SPFPPartitioner
2 from genforge.gpclassifier import gpclassifier
3 import numpy as np
4
5 # X: (N, D), y: (N,)
6 X, y = np.load("X.npy"), np.load("y.npy")
7
8 # 1) Build P disjoint (or partially overlapping) groups
9 spfp = SPFPPartitioner(n_groups=5, n_bins=100, objective="micor",
10                         perfs=0.10, remp=0.60, random_state=13, verbose=1)
11 spfp.fit(X, y)
12 groups = spfp.get_partitions() # e.g., [[2,7,5,...], [1,3,...], ...]
13
14 # 2) Hand groups to GP via userdata_pop_idx
15 params = {
16     "userdata_name": "adv_demo",
17     "userdata_xtrain": X, "userdata_ytrain": y,
18     "userdata_xval": X, "userdata_yval": y, # toy split for demo
```

```

19 "userdata_xtest": X, "userdata_ytest": y,
20 "userdata_pop_idx": groups, # <- key line
21 "runcontrol_generations": 30,
22 "runcontrol_pop_size": 64,
23 "runcontrol_useparallel": True,
24 "runcontrol_n_jobs": 8,
25 "functions_name": [["times","minus","plus","divide"]]
26 }
27
28 gp = gpclassifier.initialize(**params)
29 gp.evolve()
30 gp.report(ensemble_row=0)

```

**Notes.**

- Ensure `userdata_pop_idx` has length  $P$  (number of populations).
- If you omit `userdata_pop_idx`, each population can use all features.
- You may reuse the same SPFP partition across runs for fair comparisons.

## 10.2 Custom Primitive Function Sets

### 10.2.1 Selecting built-in tokens

Specify per-population tokens in `functions_name`:

**Listing 10.2:** Per-population primitive sets.

```

1 params["functions_name"] = [
2     ["times","minus","plus","divide"], # pop 0
3     ["times","minus","plus","divide"], # pop 1
4     ["times","minus","plus"], # pop 2
5     ["times","plus"], # pop 3
6 ]

```

The tokens `"times"`, `"minus"`, `"plus"` are always available. If your build includes protected division, `"divide"` performs safe division (avoid NaNs/inf). Additional tokens included in your local fork/build can be added the same way.

### 10.2.2 Passing custom callables

You can inject Python callables via `functions_function`. Each callable must accept `numpy` arrays and be *elementwise* safe.

**Listing 10.3:** Augmenting with custom protected ops.

```

1 import numpy as np
2
3 def safe_div(a, b, eps=1e-12):
4     return a / (np.where(np.abs(b) < eps, np.sign(b)*(eps), b))
5
6 def safe_log(a, eps=1e-12):
7     return np.log(np.where(a > eps, a, eps))
8
9 params["functions_name"] = [["times", "minus", "plus"]] # tokens
10 params["functions_function"] = [safe_div, safe_log]] # extra callables

```

#### Best practices.

- Make all custom ops *broadcasting*-compatible on `numpy` arrays.
- Guard against domain errors (e.g., negative input to `log`, zero denominator in `divide`).
- Keep the set compact; too many primitives can slow convergence.

## 10.3 Constants & Terminals

Random constants (ERCs) are governed by `const_*` keys. The exact knobs exposed by your build include (not exhaustive):

- `const_p_ERC`: probability of injecting an ERC terminal,
- `const_min`, `const_max`: numeric range (if provided),
- `const_range`, `const_about`: alternative range controls,
- `const_p_int`: bias toward integer ERCs,
- `const_num_dec_places`: decimal precision for ERCs.

**Listing 10.4:** Tuning ERC behaviour.

```

1 params.update({
2     "const_p_ERC": 0.10,
3     "const_min": -5.0, "const_max": 5.0,
4     "const_p_int": 0.20,
5     "const_num_dec_places": 3
6 })

```

Feature terminals (input variables) are wired automatically; ERCs supplement them to improve expressivity.

## 10.4 Tree & Gene Constraints

Use `gene_*` and `tree_*` to bound representation size and depth:

**Listing 10.5:** Constraining structure complexity.

```

1 params.update({
2     "gene_max_genes": [5], # per-pop gene cap
3     "tree_max_depth": [10], # per-pop tree depth cap
4     # "gene_min_genes": [...], # if exposed in your build
5     # "tree_min_depth": [...],
6 })

```

Smaller caps typically improve runtime and generalisation; increase gradually if underfitting.

## 10.5 Selection & Ensemble Mixing

Key controls reflected in your examples and sources:

- `selection_elite_fraction`: per-population elitism rate.
- `selection_elite_fraction_ensemble`: ensemble elitism rate.
- `selection_p_ensemble`: probability of participating in the ensemble.

The linear mixing head is trained with softmax-style settings (check your build for the available knobs; typical names include `softmax_*`):

**Listing 10.6:** Ensemble participation and softmax training.

```

1 params.update({
2     "selection_elite_fraction": [0.05],
3     "selection_elite_fraction_ensemble": [0.05],
4     "selection_p_ensemble": [0.50],
5
6     # Examples (adjust to your build's exact names/defaults):
7     "softmax_epochs": [50],
8     "softmax_lr": [1e-2],
9     "softmax_batch_size": [64]
10 })

```

Use `gp.report(ensemble_row=...)` to write an HTML snapshot (trees, mixing weights, metrics) for a chosen ensemble member.

## 10.6 Parallel Fitness Evaluation

GenForge supports CPU-side parallel evaluation. Typical flags (see your examples):

**Listing 10.7:** Parallel settings.

```

1 params.update({
2     "runcontrol_useparallel": True,
3     "runcontrol_n_jobs": 8, # number of worker processes/threads
4     "runcontrol_batch_job": 32 # granularity per worker (if available)
5 })

```

### Tips.

- Balance `runcontrol_batch_job` to reduce overhead on small trees.
- Fix seeds for reproducibility when running parallel (see below).

## 10.7 Saving, Logging, & Reproducibility

- Plots: enable `runcontrol_plotfitness`, `runcontrol_plotrankbest`, `runcontrol_plotrankall` and set `runcontrol_plotsavefig`.

- Plot destination: `runcontrol_plotfolder` (if exposed) controls output directory; otherwise files are saved under the run's folder.
- Snapshots: `runcontrol_savefreq` (when present) controls periodic checkpointing of run state.
- Naming: `userdata_name` tags outputs and reports.

For reproducibility, fix `random_state` (SPFP), your GP initialisation seeds (if exposed), and pin dependency versions.

## 10.8 Performance Tuning Checklist

- (a) Start small: reduce `runcontrol_generations`, `runcontrol_pop_size`, `tree_max_depth` to sanity-check the pipeline.
- (b) Use `SPFPPartitioner` to reduce the effective input dimensionality per population; tune `n_groups` and `perfs`.
- (c) Keep the primitive set lean (`functions_name`); add operators only when they demonstrably help.
- (d) Enable parallel evaluation once correctness is verified.
- (e) Inspect HTML reports to prune unhelpful primitives or raise depth caps.

## 10.9 Minimal End-to-End Template

**Listing 10.8:** Skeleton script for advanced runs.

```

1 import numpy as np
2 from genforge.spfp.spfp_partition import SPFPPartitioner
3 from genforge.gpregressor import gpregressor
4
5 # 0) Load data
6 X_train = np.load("X_train.npy")
7 y_train = np.load("y_train.npy")
8 X_val = np.load("X_val.npy")
9 y_val = np.load("y_val.npy")
10 X_test = np.load("X_test.npy")

```

```

11 y_test = np.load("y_test.npy")
12
13 # 1) Optional: SPFP partitioning
14 spfp = SPFPPartitioner(n_groups=4, n_bins=100, objective="mi",
15                         perfs=0.10, remp=0.60, random_state=19, verbose=1)
16 spfp.fit(X_train, y_train)
17 groups = spfp.get_partitions()
18
19 # 2) GP configuration
20 params = {
21     "userdata_name": "adv_reg",
22     "userdata_xtrain": X_train, "userdata_ytrain": y_train,
23     "userdata_xval": X_val, "userdata_yval": y_val,
24     "userdata_xtest": X_test, "userdata_ytest": y_test,
25     "userdata_pop_idx": groups,
26
27     "runcontrol_generations": 50,
28     "runcontrol_pop_size": 64,
29     "runcontrol_useparallel": True,
30     "runcontrol_n_jobs": 8,
31     "runcontrol_plotfitness": True,
32     "runcontrol_plotrankbest": True,
33     "runcontrol_plotsavefig": True,
34
35     "functions_name": [["times", "minus", "plus", "divide"]],
36     "gene_max_genes": [5],
37     "tree_max_depth": [10],
38
39     "selection_elite_fraction": [0.05],
40     "selection_elite_fraction_ensemble": [0.05],
41     "selection_p_ensemble": [0.50],
42
43     # softmax_* if present in your build
44 }
45
46 gp = gpregressor.initialize(**params)
47 gp.evolve()
48 gp.report(ensemble_row=0)

```

## 10.10 Troubleshooting Extensions

- **Custom op returns NaN/inf.** Add guards (e.g., `safe_div`, `safe_log`) and validate on random arrays before registering.
- **Very slow runs.** Reduce `tree_max_depth`, prune `functions_name`, and raise `runcontrol_batch_job`.
- **Ensemble not improving.** Increase `selection_p_ensemble` slightly, tune `softmax_*` epochs/lr.
- **Partitions too small/large.** Adjust `SPFPPartitioner perfs` and `n_groups`; consider `objective="micor"` to balance MI and correlation.