# GAMA

- narzędzie do automatycznego uczenia maszynowego
- potrafi automatycznie wykonać preprocessing, dobrać model oraz wykonać optymalizację hiperparametrów
- ma swoje repozytorium na Githubie
- prosty w obsłudze
- mało popularny: 62 gwiazdki na githubie, 0 cytowań

# Instalacja

Zwykła wersja

Wersja z dodatkami

```
pip install gama
```

```
pip install gama[OPTIONAL]
```

# Sposób działania

```python
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.metrics import log_loss, accuracy_score
from gama import GamaClassifier

if __name__ == "__main__":
    X, y = load_breast_cancer(return_X_y=True)
    X_train, X_test, y_train, y_test = train_test_split(
        X, y, stratify=y, random_state=0
    )

    automl = GamaClassifier(max_total_time=180, store="nothing", n_jobs=1)
    print("Starting `fit` which will take roughly 3 minutes.")
    automl.fit(X_train, y_train)

    label_predictions = automl.predict(X_test)
    probability_predictions = automl.predict_proba(X_test)

    print("accuracy:", accuracy_score(y_test, label_predictions))
    print("log loss:", log_loss(y_test, probability_predictions))
```

# Parametry

- scoring = (str, Metric, Tuple)
- regularize_length = bool
- max_pipeline_length = (int, optional)
- config = dict
- random_state = int
- verbosity = int
- search = BaseSearch
- post_processing = BasePostProcessing
- output_directory = (str, optional)
- store = str

# Parametry zasobów

- n_jobs = (int, optional)
- max_total_time = int
- max_eval _time = int
- max_memory_mb = int

# Wgranie danych z pliku

```python
from gama import GamaClassifier

if __name__ == "__main__":
    file_path = "../tests/data/breast_cancer_{}.arff"

    automl = GamaClassifier(max_total_time=180, store="nothing", n_jobs=1)
    print("Starting `fit` which take roughly 3 minutes.")
    automl.fit_from_file(file_path.format("train"))

    label_predictions = automl.predict_from_file(file_path.format("test"))
    probability_predictions = automl.predict_proba_from_file(file_path.format("test"))
```

# Preprocessing

Zmienne kategoryczne
- OneHotEncoder  (≤ 10)
- OrdinalEncoder ( > 10)
- TargetEncoder ( > 10)
- Imputowane medianą

Zmienne numeryczne
- MinMaxScaler
- MaxAbsScaler
- StandardScaler
- Normalizer
- PolynomialFeatures
- Nystroem
- RBFSampler
- PCA

# Zastosowane modele

Klasyfikacyjne
- modele bayesowskie
- regresja logistyczna
- SVM
- K najbliższych sąsiadów
- komitety modeli

Regresyjne
- modele liniowe
- SVM
- K najbliższych sąsiadów
- komitety modeli

# Algorytmy poszukiwań

- Random Search
- Asynchronous Evolutionary Algorithm (default)
- Asynchronous Successive Halving Algorithm

# Post - processing

- NoPostProcessing
- BestFitPostProcessing (default)
- EnsemblePostProcessing

# Logging

```python
import logging
import sys
from gama import GamaClassifier

gama_log = logging.getLogger('gama')
gama_log.setLevel(logging.DEBUG)

fh_log = logging.FileHandler('logfile.txt')
fh_log.setLevel(logging.DEBUG)
gama_log.addHandler(fh_log)

# The verbosity hyperparameter sets up an StreamHandler to `stdout`.
automl = GamaClassifier(max_total_time=180, verbosity=logging.DEBUG, store="nothing")
```
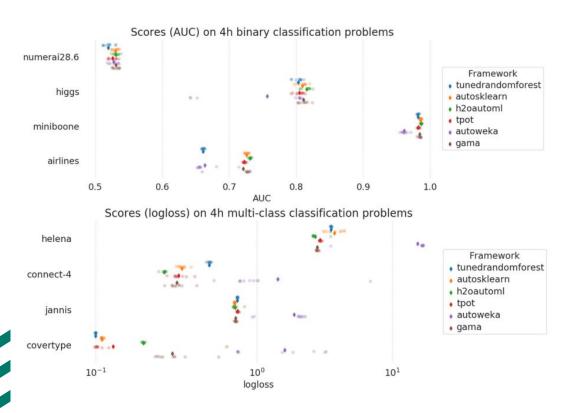
# Events

```python
from gama import GamaClassifier

def print_evaluation(evaluation):
    print(f'{evaluation.individual.pipeline_str()} was evaluated. Fitness is {evaluation.score}.')

automl = GamaClassifier()
automl.evaluation_completed(print_evaluation)
automl.fit(X, y)
```
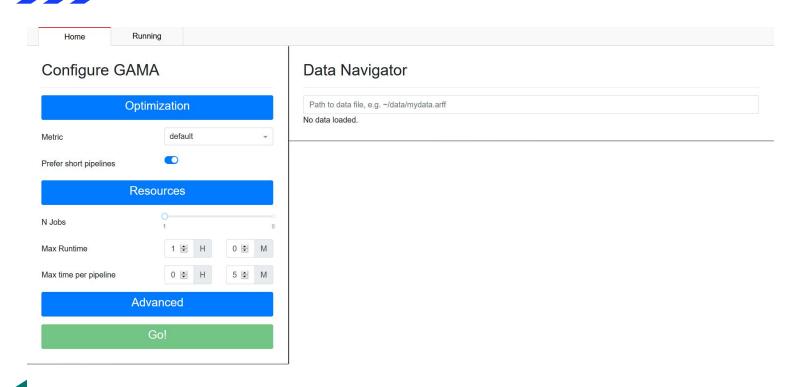
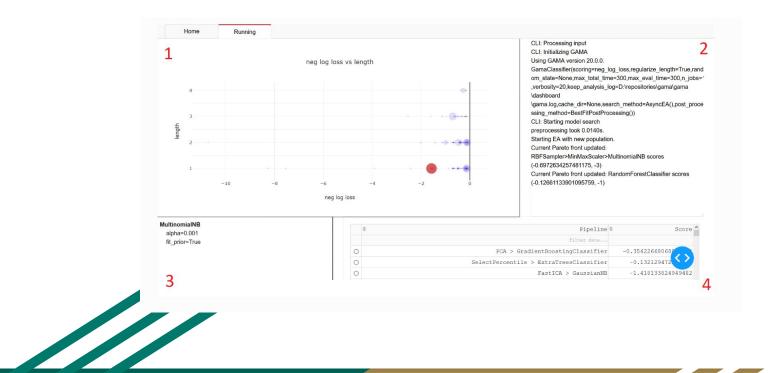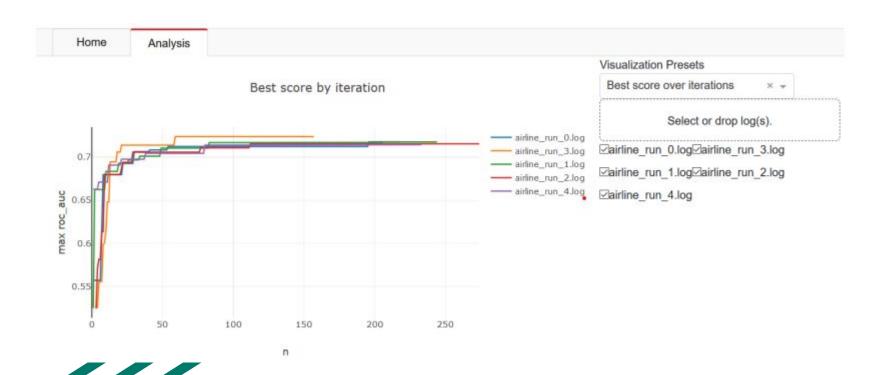# Dodawanie własnych metod wyszukiwania i post-processingu

# Benchmarki



Scores (AUC) on 4h binary classification problems

Scores (logloss) on 4h multi-class classification problems

# Dashboard

| Home | Running |
|------|---------|

## Configure GAMA

### Optimization

Metric — default ▾

Prefer short pipelines ⬤

### Resources

N Jobs — 1 ──────────── 8

Max Runtime — 1 ▲▼ H   0 ▲▼ M

Max time per pipeline — 0 ▲▼ H   5 ▲▼ M

### Advanced

### Go!

## Data Navigator

Path to data file, e.g. ~/data/mydata.arff

No data loaded.

# Running tab

# Analysis tab

# Pierwsze doświadczenia

pip install -i category_encoders==2.3.0


pip install -i Werkzeug==2.0.0

# Pierwsze pozytywne doświadczenia

```python
[5]: import logging
automl = gama.GamaClassifier(
    search=gama.search_methods.AsynchronousSuccessiveHalving(),
    post_processing=gama.postprocessing.BestFitPostProcessing(),
    n_jobs = 3,
    max_total_time=300, store="models", scoring="accuracy",
    verbosity=logging.INFO)

automl.fit(X_train, y_train)
```

```
Using GAMA version 21.0.1.
INIT:GamaClassifier(scoring=accuracy,regularize_length=True,max_pipeline_length=None,random_state=None,max_total_time=300,max_eval_time=None,n_jobs=3,max_memory_mb=None,verbosity=20,s
earch=AsynchronousSuccessiveHalving(),post_processing=BestFitPostProcessing(),output_directory=gama_7152c6ff-1365-44d6-9c0f-ed1ef0122779,store=models,goal=simplicity)
START: preprocessing default
STOP: preprocessing default after 0.0050s.
START: search AsynchronousSuccessiveHalving
ASHA start
ASHA ended due to timeout.
[2609] 3
[7830] 2
[23498] 1
Search phase evaluated 33937 individuals.
STOP: search AsynchronousSuccessiveHalving after 270.0940s.
START: postprocess BestFitPostProcessing
STOP: postprocess BestFitPostProcessing after 0.0190s.
```

```python
[6]: automl.score(X_test, y_test)
```

```
[6]: 0.9577777777777777
```

```
Data has too many features to include PolynomialFeatures
```

# Kolejne mniej pozytywne niespodzianki



```
UnboundLocalError: local variable 'highest_rung_reached' referenced before assignment
```

```
START: postprocess EnsemblePostProcessing
Not downsampling because only 1347 samples were stored.
Error during auto ensemble: division by zero
Traceback (most recent call last):
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\site-packages\gama\postprocessing\ensemble.py", line 524, in build_fit_ensemble
    ensemble.build_initial_ensemble(10)
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\site-packages\gama\postprocessing\ensemble.py", line 265, in build_initial_ensemble
    self._ensemble_validation_score()
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\site-packages\gama\postprocessing\ensemble.py", line 444, in _ensemble_validation_score
    prediction_to_validate = self._averaged_validation_predictions()
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\site-packages\gama\postprocessing\ensemble.py", line 240, in _averaged_validation_predictions
    return weighted_sum_predictions / self._total_model_weights()
ZeroDivisionError: division by zero
STOP: postprocess EnsemblePostProcessing after 0.0080s.
```
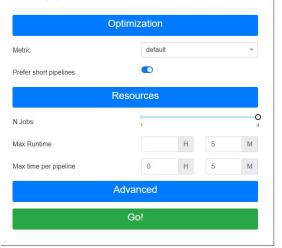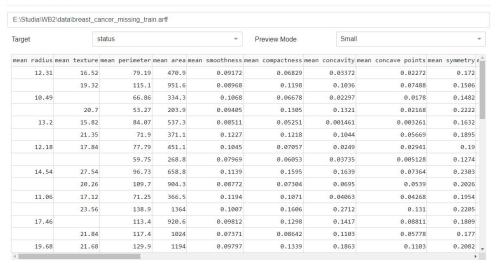
# Jeszcze mniej pozytywne niespodzianki

```
Traceback (most recent call last):
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\multiprocessing\queues.py", line 241, in _feed
    send_bytes(obj)
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\multiprocessing\connection.py", line 200, in send_bytes
    self._send_bytes(m[offset:offset + size])
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\multiprocessing\connection.py", line 280, in _send_bytes
    ov, err = _winapi.WriteFile(self._handle, buf, overlapped=True)
BrokenPipeError: [WinError 232] Trwa zamykanie potoku
Traceback (most recent call last):
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\multiprocessing\queues.py", line 241, in _feed
    send_bytes(obj)
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\multiprocessing\connection.py", line 200, in send_bytes
    self._send_bytes(m[offset:offset + size])
  File "C:\Users\PC-Komputer\Anaconda3\envs\gama\lib\multiprocessing\connection.py", line 280, in _send_bytes
    ov, err = _winapi.WriteFile(self._handle, buf, overlapped=True)
BrokenPipeError: [WinError 232] Trwa zamykanie potoku
```
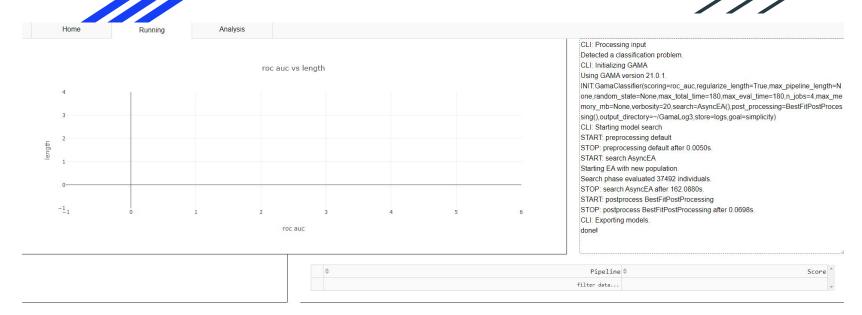
# Zadziwiające pozytywne doświadczenia

# Które okazało się złudne

# Dziękujemy za obejrzenie prezentacji

Tomasz Siudalski, Grzegorz Zbrzeżny, Piotr Marciniak

# Źródła:

- GAMA: a General Automated Machine learning Assistant
- https://github.com/openml-labs/gama
- https://openml-labs.github.io/gama/master