

# PYTHON CHO KHOA HỌC DỮ LIỆU

Giảng viên: Hà Minh Tuấn

## THAM SỐ MẶC ĐỊNH CỦA ModelTrainer

Nhóm thực hiện:

Dặng Dĩnh Đoàn 23280046  
Phạm Thanh Uy 23280097

TP. Hồ Chí Minh, tháng 12 năm 2025

# Mục lục

<b>1</b>	<b>Giới thiệu</b>	<b>2</b>
<b>2</b>	<b>Nhóm cấu hình chính – TrainerConfig</b>	<b>2</b>
2.1	Target và random state . . . . .	2
2.2	Mục tiêu tối ưu và metrics . . . . .	2
2.3	Cross-validation & search . . . . .	3
2.4	Imbalance & danh sách model . . . . .	3
2.5	Thư mục lưu model . . . . .	3
<b>3</b>	<b>Nhóm “nút vặn” ẩn trong method</b>	<b>4</b>
3.1	run_training(...) . . . . .	4
3.2	_get_model_and_param_grid(...) . . . . .	4
3.3	evaluate_on_test(...) . . . . .	5
3.4	compute_shap_values(...) . . . . .	5
3.5	Lưu / load model . . . . .	5
3.6	RESULTS_DIR . . . . .	5

# 1 Giới thiệu

Đây là các tham số mặc định của class ModelTrainer, khi User không truyền vào, sẽ tự động lấy các tham số này.

## 2 Nhóm cấu hình chính – TrainerConfig

```
@dataclass
class TrainerConfig:
    target_col: str = "Outcome"
    random_state: int = 42
    scoring_primary: str = "f1"
    scoring_other: List[str] = ["roc_auc", "accuracy", "precision", "recall"]
    cv_splits: int = 5
    use_randomized_search: bool = True
    n_iter_random_search: int = 20
    use_smote: bool = True
    model_names: List[str] = ["log_reg", "random_forest"]
    model_output_dir: Optional[Union[str, Path]] = None
```

### 2.1 Target và random state

- **target\_col**
  - Hiện tại: "Outcome"
  - Có thể đổi theo dataset khác: "Label", "HasDiabetes", "Target", ...
- **random\_state**
  - Hiện tại: 42
  - Có thể đổi sang bất kỳ số nguyên: 0, 123, 2024, ...
  - Nên cố định để kết quả reproducible

### 2.2 Mục tiêu tối ưu và metrics

- **scoring\_primary**
  - Hiện tại: "f1"
  - Các option điển hình của scikit-learn:
    - \* "roc\_auc", "accuracy", "precision", "recall", "balanced\_accuracy"
    - \* Hoặc custom scorer nếu tự định nghĩa
- **scoring\_other**
  - Hiện tại: ["roc\_auc", "accuracy", "precision", "recall"]
  - Có thể thêm/bớt metric phụ:
    - \* Bớt: ["roc\_auc", "accuracy"]

- \* Thêm: ["roc\_auc", "accuracy", "precision", "recall", "balanced\_accuracy"]
- Lưu ý: phải là tên metric mà sklearn hỗ trợ

## 2.3 Cross-validation & search

- **cv\_splits**
  - Hiện tại: 5
  - Có thể đổi: 3, 4 (nhanh hơn), 10 ( ổn định hơn nhưng chậm)
- **use\_randomized\_search**
  - True: dùng RandomizedSearchCV
  - False: dùng GridSearchCV
- **n\_iter\_random\_search**
  - Hiện tại: 20
  - Chỉ áp dụng khi use\_randomized\_search=True
  - Có thể giảm: 10 (nhanh hơn), tăng: 50 (tìm kỹ hơn, chậm hơn)

## 2.4 Imbalance & danh sách model

- **use\_smote**
  - Hiện tại: True
  - True: dùng SMOTE để cân bằng nhãn
  - False: không dùng SMOTE, rely vào class\_weight='balanced'
- **model\_names**
  - Hiện tại: ["log\_reg", "random\_forest"]
  - Có thể chỉnh:
    - \* Chỉ Logistic: ["log\_reg"]
    - \* Chỉ Random Forest: ["random\_forest"]
    - \* Mở rộng: ["log\_reg", "random\_forest", "xgboost"], ["log\_reg", "svm"], ...

## 2.5 Thư mục lưu model

- **model\_output\_dir**
  - Hiện tại: None → không lưu model
  - Có thể đặt:
    - \* "models" → <project\_root>/models
    - \* "artifacts/models" → <project\_root>/artifacts/models

\* Đường dẫn tuyệt đối: r"E:\= "

### 3 Nhóm “nút văn” ẩn trong method

#### 3.1 run\_training(...)

- X\_test, y\_test
  - Nếu truyền: tính metrics test, experiments.csv, best\_model.txt
  - Nếu None: chỉ CV trên train, không có test metrics
- n\_jobs
  - Hiện tại: n\_jobs=1
  - Có thể đổi: -1 để dùng hết core
- verbose
  - Hiện tại: 1
  - Có thể tăng: 2 để log chi tiết

#### 3.2 \_get\_model\_and\_param\_grid(...)

- Logistic Regression
  - solver: "liblinear" (hiện dùng), "lbfgs", "saga"
  - max\_iter: 1000 (có thể tăng 2000)
  - class\_weight: "balanced" / None
  - param\_grid:
    - \* C: [0.01, 0.1, 1.0, 10.0] (có thể thêm log-scale)
    - \* penalty: ["l1", "l2"] (phụ thuộc solver)
- Random Forest
  - n\_estimators: 200 (tăng 300, 500 nếu muốn)
  - class\_weight: "balanced" / None
  - n\_jobs: 1 (có thể -1)
  - param\_grid:
    - \* n\_estimators: [100, 200, 300]
    - \* max\_depth: [None, 5, 10, 20]
    - \* min\_samples\_split: [2, 5]
- SMOTE
  - True: Pipeline(SMOTE + model), prefix "model\_\_" cho param\_grid

- False: model gốc, param\_grid không có prefix

### 3.3 evaluate\_on\_test(...)

- Metrics hiện tại: accuracy, precision, recall, f1, roc\_auc, tn, fp, fn, tp
- Có thể thêm: specificity = tn / (tn+fp), balanced\_accuracy, ...
- zero\_division: 0 (có thể đổi 1 hoặc "warn")

### 3.4 compute\_shap\_values(...)

- max\_samples
  - Default: 200
  - Khi gọi: trainer.compute\_shap\_values("log\_reg", X\_train, max\_samples=500)
  - Ít → nhanh hơn, nhiều → ổn định hơn nhưng chậm
- Lớp multi-class
  - Nếu sv.ndim=3 và sv.shape[1]=2 → lấy lớp positive (index 1)
  - Ngược lại → average theo axis=1

### 3.5 Lưu / load model

- \_maybe\_save\_model(...)
  - Tên file: model\_output\_dir / f"{{model\_name}}\_best.joblib"
  - Có thể đổi pattern: f"{{model\_name}}\_best\_{scoring\_primary}.joblib"
- load\_saved\_model(model\_name, path=None)
  - Truyền path cụ thể nếu muốn load khác
  - path=None → dùng <model\_output\_dir>/<model\_name>\_best.joblib

### 3.6 RESULTS\_DIR

- Thư mục results dùng trong:
  - experiments.csv, best\_model.txt
- Có thể chỉnh trong src/utils/config.py:
  - Đổi RESULTS\_DIR sang Path("outputs"), Path("reports"), ...