Recommended Environment:

Run this notebook in the model eval suite Conda environment for best results.

See setup instructions in the Usage Guide.

If you're running this outside Conda, you can install dependencies manually: Uncomment the line below to install from the root requirements file.

!pip install -r ../../requirements.txt

Model Evaluation Suite Demo Notebook

This notebook demonstrates how to use the Model Evaluation Suite to:

- Prepare and validate input data.
- Run a full modeling pipeline using YAML configuration files.
- · Log models and artifacts with MLflow.
- Evaluate a production candidate against a holdout dataset.
- Optionally compare against a baseline model for performance drift or uplift.
- Project Structure
- Workflow Overview

Configuration Setup

This notebook is driven by modular YAML configuration files, which serve as the central control system for the evaluation suite.

click here to expand section

🔧 Custom Feature Engineering

This suite supports plug-and-play custom transformers via the feature_engineering block in your YAML.

click here to expand



😉 Dashboard Guidence



Pre-Model Diagnostics Dashboard

This optional module runs before any model training or validation occurs. It provides key insights into the integrity and statistical structure of your input data. It is driven by the pre_model_diagnostics block in your YAML and is best used in notebook workflows.

click here to expand section

Model Evaluation Dashboard

This dashboard provides an interactive summary of the model trained in the experimental run. It visualizes performance on the test set and includes explainability tools to support model diagnostics and stakeholder communication.

- ► Summary
- ▶ Details
- ► Importance
- ► Explainability
- ▶ Plotviewers
- ▶ Metadata
- ▶ Alerts

Core Imports

You can access the main runners directly from the package thanks to a clean interface exposed via __init__.py . These entrypoints allow you to run each stage of the pipeline from a single import.

```
In [1]: from model_eval_suite import run_experiment, validate_champion, prep_data
```

```
In [2]: import os
# Change the working directory to the project root
    os.chdir('...')
```

Data Preparation

If you're starting from raw CSVs, you can use the suite's built-in preprocessing tool, data_prep.py to split the data into training, testing, and holdout sets.

Skip this if you've already created your train.csv, test.csv, and holdout.csv.

```
In [3]: prep_data(config_path="config/data_prep.yaml")
Loading raw data from: data/input_data/salifort_50k.csv
```

Performing train/test split on development data...

☑ Train data saved to: data/dev_data/train_data.csv (30000 rows)

Test data saved to: data/dev_data/test_data.csv (10000 rows)

☑ Holdout data saved to: data/holdout data/holdout data.csv (10000 rows)

Model Experiment Runs (Demo)

This demo walks through multiple model runs using the salifort 50k dataset. Although the dataset is optimized for classification tasks, it is used for both classification and regression pipelines to demonstrate flexibility and YAML-driven control.

We run the following models using the evaluation suite:

Classifier Models

• Gaussian Naive Bayes

Config: config/classifier/guas_nb.yaml

Logistic Regression

Config: config/classifier/logreg.yaml

XGBoost Classifier

Config: config/classifier/xgboost.yaml

Regressor Models

• Linear Regression

Config: config/regressor/linreg.yaml

XGBoost Regressor

Config: config/regressor/xgboost reg.yaml

Each model triggers:

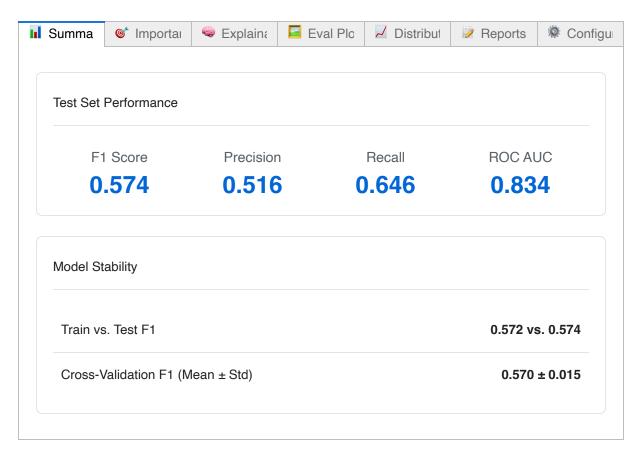
- An optional **Pre-Model Diagnostics Dashboard** (if enabled in YAML)
- A complete Evaluation Dashboard with explainability, distributions, and exportable artifacts

At the end of the demo, we use the **champion validation** system to validate and crown the two XGBoost models — one for classification, one for regression.

All behavior is controlled by the YAML configs. See the config/ directory or config resources/config.zip for template downloads.

In [4]: # ======= Naive - Baves ======= run_experiment(user_config_path="config/classifier/guas_nb.yaml") Registered model 'gnb_demo_01' already exists. Creating a new version of thi s model... Created version '4' of model 'gnb_demo_01'. SHAP explainer creation failed: The passed model is not callable and cann ot be analyzed directly with the given masker! Model: GaussianNB() Run complete: `gnb_demo_01` Artifacts saved to: - Plots: exports/plots/gnb_demo_01 - Reports: exports/reports/gnb_demo_01 MLflow model: `qnb demo 01` --- Rendering Dashboards ---☐ Pre-Modeling Diagnostics Overview Missingness Collinearity Distributions Eval Plots Skewed Features **Outlier Counts** skewness outlier count time_spend_company ||1.924189 time_spend_company | 4420 Work_accident 1.922905 Work_accident 7709 promotion_last_5years 6.426524 promotion last 5years 1400

Run: gnb_demo_01 | **F1 (Test):** 0.574



SHAP Error Handling

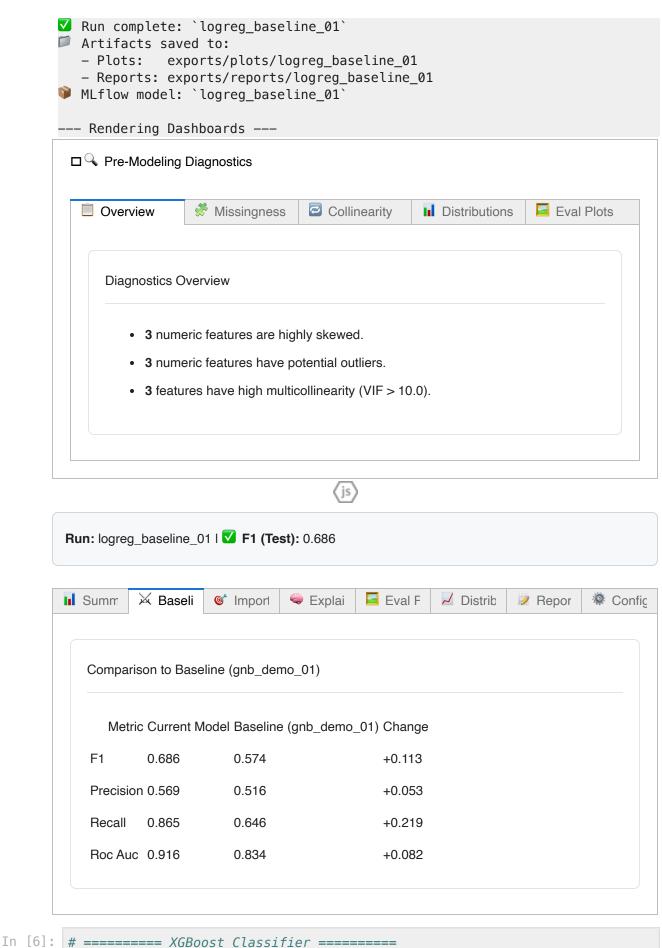
▼ click to expand section

Some models — such as SVC, SVR, or GaussianNB — do not expose traditional feature importance attributes or are incompatible with SHAP explainability tools.

This suite handles such situations gracefully:

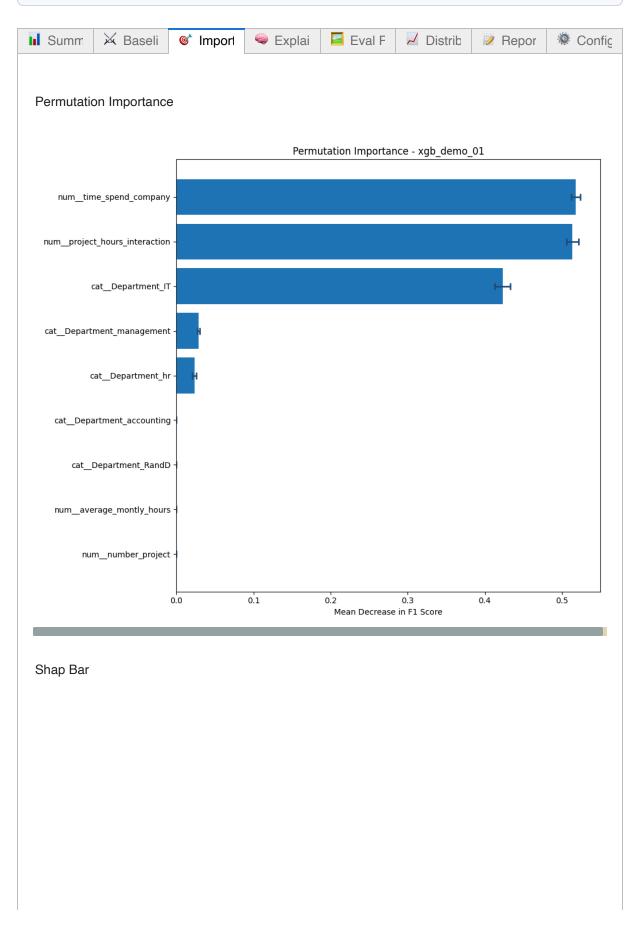
- The SHAP tab will be skipped silently if no compatible features are found.
- A warning will be logged (but not treated as a failure).
- All other evaluation plots and metrics will still render normally.

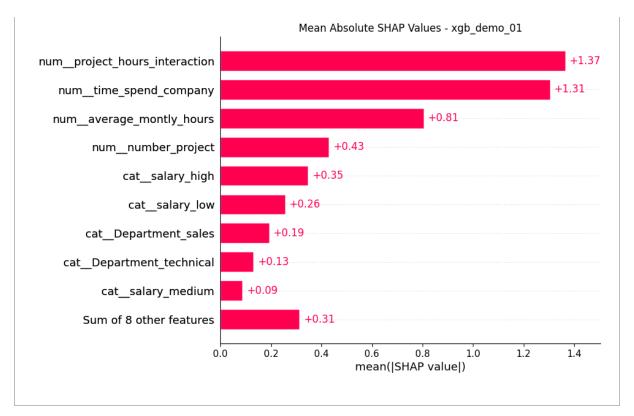
This ensures that the workflow remains robust even for models with limited explainability tooling.



 $file: ///Users/garrettschumacher/Documents/git_repos/model_evaluation_suite/notebooks/demo.html$

run_experiment(user_config_path="config/classifier/xgboost.yaml") Fitting 5 folds for each of 108 candidates, totalling 540 fits Registered model 'xgb_demo_01' already exists. Creating a new version of thi s model... Created version '4' of model 'xgb_demo_01'. Run complete: `xgb_demo_01` Artifacts saved to: - Plots: exports/plots/xgb_demo_01 - Reports: exports/reports/xgb_demo_01 MLflow model: `xgb_demo_01` --- Rendering Dashboards ---☐ Pre-Modeling Diagnostics Overview Missingness Eval Plots Collinearity Distributions VIF Scores VIF feature 0 last_evaluation 18.815598 average_montly_hours 18.091514 2 number_project 13.411024 satisfaction_level 6.689082 time_spend_company 6.150695 Work_accident 1.184356 promotion_last_5years 1.030955



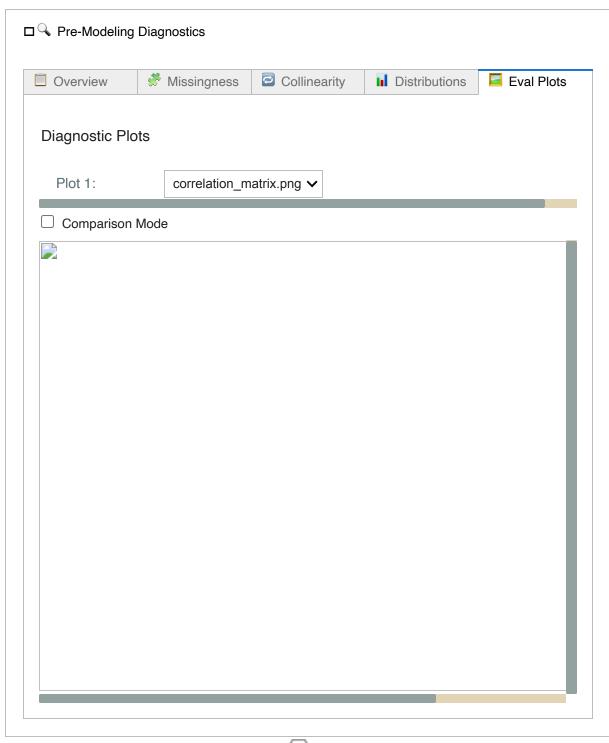


Cross-Validation Insight

▼ click here to expand

If hyperparameter tuning via cross-validation is enabled in the config (hyperparameter_tuning.run: true), the dashboard will include an additional boxplot in the **Summary** tab.

This plot visualizes the distribution of CV scores across folds for the best-performing parameter set, offering a quick diagnostic of stability and performance variance.



Run: linreg_demo_01 | R-squared: 0.331

Audit Alerts:

• 🚣 Weak Fit: R-squared is low (0.331), indicating poor model fit.

| Summa 🏻 🍑 Im | | Explaina | Eval Plo | ✓ Distrib | out Pepo | orts 🦃 Co |
|--------------------------|---------|------------|-------------|-------------|-----------|-----------|
| Statistical Summ | ary | | | | | |
| | | | | | | |
| | OLS | Regressio | n Results | | | |
| | ===== | | ====== | ====== | ======= | |
| Dep. Variable | | | -squared: | | 0.336 | |
| Model: | | | dj. R-squar | ed: | 0.336 | |
| Method: | Lea | st Squares | s F-statist | ic: | 950.0 | |
| Date: | Tue, 05 | 5 Aug 2025 | 5 Prob (F- | statistic): | 0.0 | 00 |
| Time: | 20 |):52:48 L | og-Likeliho | od: | -7042.4 | |
| No. Observati | ons: | 3000 | 0 AIC: | | 1.412e+0 | 04 |
| Df Residuals: | | 29983 | BIC: | | 1.426e+04 | |
| Df Model: | | 16 | | | | |
| Covariance Ty | /pe: | nonrobu | ust | | | |
| | | | | | = | |
| | | coef sto | l err t | P>Itl | [0.025 | 0.975] |
| const | | -9.38e+09 | 8.47e+10 | 0.111 | 0.912 | -1.75e+11 |
| 1.57e+11 | | | | | | |
| numnumbe | | - | 0.7656 | 0.007 -1 | 11.520 0 | .000 -0 |
| | | , hours | 0 5000 | 0.006 | 102 017 | 0.000 |
| numaverag -0.593 -0.5 | | _110018 | -0.3622 | 0.006 | -103.017 | 0.000 |

| numtime_spend_company 0.040 0.048 | 0.0440 0.002 24.221 0.000 |) |
|---|---|-----|
| | on 1.1377 0.010 114.724 0.000 | |
| catDepartment_IT - 99e+10 2.67e+10 | 1.596e+09 1.44e+10 -0.111 0.912 | -2. |
| catDepartment_RandD -2.99e+10 2.67e+10 | -1.596e+09 1.44e+10 -0.111 0.91 | 2 |
| | -1.596e+09 1.44e+10 -0.111 0.9 | 12 |
| | 1.596e+09 1.44e+10 -0.111 0.912 | -2. |
| | nt -1.596e+09 1.44e+10 -0.111 0 | .91 |
| catDepartment_marketing -2.99e+10 2.67e+10 | -1.596e+09 1.44e+10 -0.111 0.9 ⁻¹ | 12 |
| catDepartment_product_mi | ng -1.596e+09 1.44e+10 -0.111 0. | 91 |
| catDepartment_sales 2.99e+10 2.67e+10 | -1.596e+09 1.44e+10 -0.111 0.912 | - |
| catDepartment_support 2.99e+10 2.67e+10 | -1.596e+09 1.44e+10 -0.111 0.912 | 2 - |
| | -1.596e+09 1.44e+10 -0.111 0.912 | 2 - |
| | 098e+10 9.92e+10 0.111 0.912 -1 | .83 |
| | 98e+10 9.92e+10 0.111 0.912 -1. | .83 |
| | 1.098e+10 9.92e+10 0.111 0.912 | -1. |
| | | === |
| | 08 Durbin-Watson: 2.016 000 Jarque-Bera (JB): 9069.694 | |

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 3.18e-26. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

Automated Alert Auditing

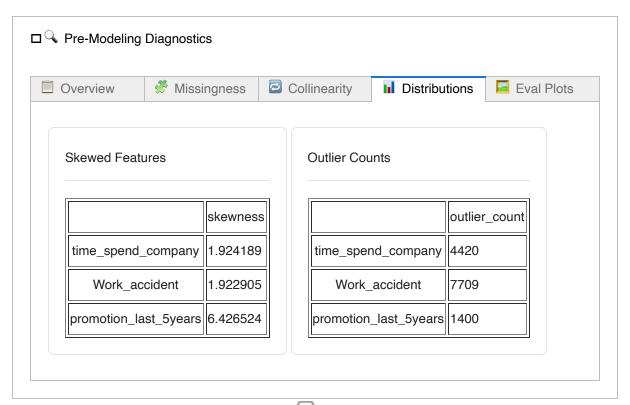
▼ click here to expand section

The validation dashboard includes an **Alerts** tab that surfaces automated audit checks on your model's performance.

These alerts are designed to flag potential concerns such as:

- Very low precision or recall
- High class imbalance
- Overfitting indicators (e.g., large delta between train/test scores)
- Underwhelming performance against a baseline (if provided)

This system provides a lightweight, interpretable review of model quality without requiring custom code or manual thresholding.



Run: xgbreg_demo_01 | R-squared: 0.854

```
Config
                     o Import
                                Explai
                                          Eval F
                                                     Distrib
Summ

→ Baseli

                                                                Repor
 "run_id": "xgbreg_demo_01",
 "task_type": "regression",
 "pre_model_diagnostics": {
   "run": true,
   "export_reports": true,
   "export_html_report": true,
   "vif_threshold": 10.0,
   "skewness threshold": 0.75
 },
  "notebook mode": true,
  "logging": "auto",
  "paths": {
   "input_data": "data/input_data/salifort_50k.csv",
   "reports_dir": "exports/reports",
   "plots dir": "exports/plots",
   "model_export_dir": "models",
   "metrics_log": "exports/reports/model_metrics_log.csv",
   "log_dir": "logs",
   "train_data_path": "data/dev_data/train_data.csv",
   "test_data_path": "data/dev_data/test_data.csv"
 },
  "modeling": {
   "target_column": "left",
   "test size": 0.2,
   "feature_engineering": {
    "run": true,
    "module": "model_eval_suite.modeling.feature_engineering",
    "class_name": "HREngineer"
```

```
"hyperparameter_tuning": {
 "run": true,
 "param_grid": {
  "estimator__n_estimators": [
    50,
   100,
   150
  ],
  "estimator__max_depth": [
    3,
    5,
    7
  ],
  "estimator__learning_rate": [
   0.01,
   0.1,
    0.2
  ],
  "estimator__subsample": [
    0.8,
    1.0
  "estimator__colsample_bytree": [
    0.8,
    1.0
 "cv_folds": 5,
 "scoring": "r2",
 "verbose": 1
},
"pipeline_factory": {
 "name": "XGBRegressor",
```

```
"registered_name": "xgbreg_demo_01",
  "numeric_features": [
   "number_project",
   "average_montly_hours",
   "time_spend_company",
   "project_hours_interaction"
  1,
  "categorical_features": [
   "Department",
   "salary"
  ],
  "params": {
   "n_estimators": 100,
   "max_depth": 5,
   "learning_rate": 0.1,
   "objective": "reg:squarederror",
   "random state": 42
},
"evaluation": {
 "run": true,
 "export_xlsx_summary": false,
 "export_html_dashboard": true,
 "compare_to_baseline": "linreg_demo_01",
 "plots": {
  "confusion_matrix": {
   "save": true
  },
  "roc curve": {
   "save": true
  },
  "pr_curve": {
   "save": true
```

```
"learning_curve": {
 "save": true
},
"calibration_curve": {
 "save": true
},
"threshold_plot": {
 "save": true
},
"permutation_importance": {
 "save": true
},
"lift_curve": {
 "save": true
},
"cumulative_gain": {
 "save": true
},
"shap_summary": {
 "save": true
},
"feature_distributions": {
 "save": true
},
"predicted_vs_actual": {
 "save": true
},
"residuals_plot": {
 "save": true
},
"residuals_histogram": {
 "save": true
},
```

```
"qq_plot": {
  "save": true
},
"explainability": {
 "run_shap": true,
 "shap_sample_size": 2000
},
"audits": {
 "r_squared_threshold": 0.6,
 "rmse threshold": 50000.0,
 "accuracy_threshold": 0.8,
 "f1 score threshold": 0.75
```

Y Champion Model Validation

This section evaluates a registered MLflow model (your champion) against a holdout dataset using a dedicated validation YAML configuration.

Key Features

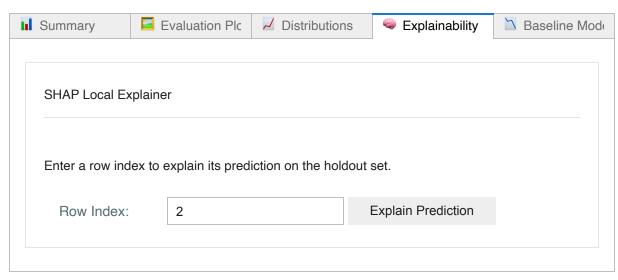
- Uses its own standalone YAML file (separate from training experiments)
- Accepts an optional **baseline model** for drift detection or performance benchmarking
- Automatically generates:
 - Confidence interval plots (if applicable)
 - Baseline comparison deltas (if a baseline model is provided)
 - Alert audits for performance degradation or instability
- Produces a complete interactive dashboard with:
 - Summary metrics and cross-validation visualizations
 - Explainability and feature importance plots
 - Distribution visualizations for target and predictions
 - Full configuration and environment metadata
- Tags the evaluated model in the MLflow Registry using your specified production_tag

This workflow is ideal for pre-deployment validation, regression testing, and model promotion decisions.

Validation Configurations Used in This Demo

- config/xgb_validation.yaml XGBoost classifier
- config/xgb_reg_validation.yaml XGBoost regressor

Champion Model Validation: xgb_demo_01_production_validation



```
Detected task type: regression

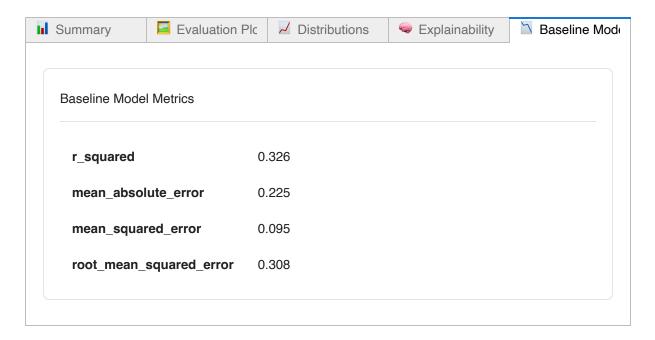
Generating final assessment plots...

Prediction interval plot is only available for ensemble models like Rando mForestRegressor.

Exporting validation artifacts...

Tagging model version with status: 'Production-Candidate'
--- Validation Complete for xgbreg_demo_01 v1 ---
--- Rendering Validation Dashboard ---
```

Champion Model Validation: xgbreg_demo_01_production_validation



☑ Wrap-Up and Next Steps

You've now run multiple models through the full suite — from preprocessing and diagnostics to evaluation and champion validation.

This notebook demonstrates the flexibility of the system, including:

- YAML-driven configuration at every stage
- Reusable pipelines for both classification and regression tasks
- Support for custom feature engineering and hyperparameter tuning
- Interactive dashboards for diagnostics and final reporting
- MLflow integration for model tracking and registry updates

Next Steps

- **Test additional models** by duplicating a config YAML.
- Customize features using your own transformers or fe_config modules.

• Enable advanced diagnostics, SHAP, and permutation importance as needed.

• Package and deploy validated models via the MLflow registry.

For more examples and config templates, explore the config_resources/ folder or the full README.md.

Questions or suggestions? Feel free to submit an issue or feature request in the repository.