

A Comprehensive Field Guide for Data Scientists

February 2026

Contents

Preface	3
1 Cross-Validation: Honest Performance Estimation	4
1.1 The Core Problem: Why a Single Train–Test Split Is Not Enough	4
1.2 K-Fold Cross-Validation	4
1.2.1 How It Works	4
1.2.2 Choosing k	4
1.2.3 Python Example: K-Fold and Stratified K-Fold	5
1.3.1 Walk-Forward Validation (Expanding Window)	7
2 Hyperparameter Tuning: Systematic Model Configuration	10
2.1 What Are Hyperparameters?	10
2.2 Grid Search	10
2.3 Random Search	12
2.4 Bayesian Optimisation	13
2.4.1 Choosing a Tuning Strategy	15
3 Metrics for Imbalanced Datasets	17
3.1 Why Accuracy Fails	17
3.2 The Confusion Matrix: A Foundation	17
3.3 Core Metrics	17
3.3.1 Reading Precision-Recall Curves	21
3.3.2 Threshold Optimisation	21
4 Model Interpretability	22
4.1 Why Interpretability Is Not Optional	22
4.2 SHAP: SHapley Additive exPlanations	22
4.2.1 The Theory	22
4.2.2 Reading SHAP Output	25
4.3 LIME: Local Interpretable Model-agnostic Explanations	25
4.3.1 SHAP vs LIME: When to Use Each	28
4.4 Partial Dependence Plots	28
5 Putting It All Together: End-to-End Evaluation Workflow	32
6 The Evaluation Checklist	36
7 Regression Metrics: Measuring Continuous Predictions	38
7.1 The Problem with a Single Number	38
7.2 Core Regression Metrics	38
7.3 Residual Analysis: Beyond a Single Number	39
8 Model Calibration: When Probabilities Must Mean Something	44
8.1 The Calibration Problem	44
8.2 The Reliability Diagram	44
8.3 The Brier Score	44

9	Statistical Significance Testing for Model Comparison	48
9.1	The Core Question: Is Model A Actually Better?	48
9.2	McNemar's Test: Comparing Classifiers on the Same Test Set	48
9.3	The Corrected Resampled t-Test	48
10	Bias & Fairness Evaluation	52
10.1	Why Fairness Is a Measurement Problem	52
10.2	Core Fairness Definitions	52
11	Production Monitoring & Drift Detection	57
11.1	The Life Cycle Does Not End at Deployment	57
11.2	Population Stability Index (PSI)	57
12	Learning Curves & Diagnosing Underfitting vs Overfitting	62
12.1	The Bias-Variance Trade-off in Practice	62
12.2	Reading Learning Curves	62
12.3	The Validation Curve: Finding the Complexity Sweet Spot	65
13	Multiclass & Multi-label Metrics	66
13.1	Beyond Binary Classification	66
13.2	Averaging Strategies	66
13.3	Cohen's Kappa: Agreement Beyond Chance	66
13.4	The Normalised Confusion Matrix as a Debugging Tool	69
	Appendix: Quick Reference	72

Preface: Why Evaluation Is the Most Important Skill

Building a machine learning model that produces numbers is easy. Building one you can *trust*, defend in a board meeting, and rely on to drive real decisions is an entirely different discipline. That discipline is **Model Evaluation & Validation**.

Many data science projects fail not because the algorithm was wrong, but because the practitioner optimised for the wrong metric, let data from the future leak into training, reported accuracy on a wildly imbalanced test set, or could not explain to a stakeholder why the model made a particular decision. This tutorial is designed to close those gaps.

The guide is structured around four pillars, each building on the previous one:

1. **Cross-Validation** — learning how to estimate generalisation performance honestly, and how to avoid the most common pitfalls like data leakage and overfitting to a validation fold.
2. **Hyperparameter Tuning** — systematically searching the space of model configurations rather than guessing, using methods that scale from exhaustive search to Bayesian optimisation.
3. **Metrics for Imbalanced Datasets** — understanding why accuracy is misleading when classes are unequal and which metrics actually reflect business cost.
4. **Model Interpretability** — opening the black box with SHAP, LIME, and partial dependence plots so that the model can be scrutinised, debugged, and explained.

Throughout, every concept is grounded in runnable Python examples using mainstream libraries (`scikit-learn`, `shap`, `lime`, `scikit-optimize`). Code is annotated line by line wherever a concept needs emphasis.

1 Cross-Validation: Honest Performance Estimation

1.1 The Core Problem: Why a Single Train–Test Split Is Not Enough

Suppose you train a classifier and measure its accuracy on a held-out test set. You get 87%. Should you celebrate? Not yet. That number reflects one particular random partition of your data. Had you used a different random seed, you might have gotten 82% or 91%. You have measured the performance of *one* model on *one* subset of your data, not the expected performance of your modelling *approach* on unseen data.

Cross-validation resolves this by evaluating the model on multiple, non-overlapping validation subsets and averaging the results. The variance of those scores tells you how stable your model is; the mean tells you the unbiased estimate of generalisation performance.

Just as a model can overfit the training data, an *evaluation procedure* can overfit the validation data. Every time you make a modelling decision based on a fixed validation set, you implicitly “train” on that set. Cross-validation delays this leakage by rotating which data is held out.

1.2 K-Fold Cross-Validation

1.2.1 How It Works

K-Fold partitions the dataset into k equally sized ***folds***. The model is trained k times; each time, one fold serves as the validation set and the remaining $k - 1$ folds form the training set. The final performance estimate is the mean (and standard deviation) across all k scores.

Fold 1	VAL	TRAIN	TRAIN	TRAIN	TRAIN
Fold 2	TRAIN	VAL	TRAIN	TRAIN	TRAIN
Fold 3	TRAIN	TRAIN	VAL	TRAIN	TRAIN
Fold 4	TRAIN	TRAIN	TRAIN	VAL	TRAIN
Fold 5	⇒ TRAIN	TRAIN	TRAIN	TRAIN	VAL

Figure 1: 5-Fold Cross-Validation: the teal fold rotates through each iteration.

1.2.2 Choosing k

- **$k = 5$ or 10** is the standard choice backed by empirical research (Kohavi 1995). It balances bias and variance in the estimator.
- **$k = n$** (Leave-One-Out CV) is nearly unbiased but has very high variance and is computationally expensive.
- **Stratified K-Fold** preserves the class distribution in each fold — always use this for classification.

1.2.3 Python Example: K-Fold and Stratified K-Fold

K-Fold Cross-Validation — Full Walkthrough

```

1  import numpy as np
2  import pandas as pd
3  from sklearn.datasets import make_classification
4  from sklearn.ensemble import RandomForestClassifier
5  from sklearn.model_selection import (
6      KFold, StratifiedKFold, cross_val_score, cross_validate
7  )
8  from sklearn.preprocessing import StandardScaler
9  from sklearn.pipeline import Pipeline
10
11  # -----
12  # 1. Synthetic dataset: 1000 samples, 20 features, binary target
13  #    weights=(0.7, 0.3) means 70% class 0, 30% class 1 mild
14  #    imbalance
15  # -----
16  X, y = make_classification(
17      n_samples=1000,
18      n_features=20,
19      n_informative=10,
20      n_redundant=4,
21      weights=[0.7, 0.3],
22      random_state=42
23  )
24  # -----
25  # 2. Build a Pipeline
26  #    Always wrap preprocessing inside the pipeline so that the
27  #    scaler
28  #    is fit ONLY on training folds never on the validation fold.
29  #    This is the single most common source of data leakage.
30  # -----
31  pipeline = Pipeline([
32      ('scaler', StandardScaler()),          # zero-mean, unit-
33      ('clf', RandomForestClassifier(        variance
34          n_estimators=100, random_state=42
35      ))
36  ])
37  # -----
38  # 3. Plain K-Fold (ignores class distribution in each fold)
39  # -----
40  kf = KFold(n_splits=5, shuffle=True, random_state=42)
41
42  scores_kfold = cross_val_score(
43      pipeline, X, y,
44      cv=kf,
45      scoring='roc_auc',          # area under ROC curve: threshold-
46      n_jobs=-1                  independent
47  )                                # use all CPU cores
48
49  print("=== Standard K-Fold ===")
50  print(f"Fold AUC scores : {np.round(scores_kfold, 4)}")

```

```

51 print(f"Mean AUC          : {scores_kfold.mean():.4f}")
52 print(f"Std AUC          : {scores_kfold.std():.4f}")
53
54 # -----
55 # 4. Stratified K-Fold (preserves class ratio in every fold)
56 #   For classification, ALWAYS prefer StratifiedKFold.
57 # -----
58 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
59
60 scores_stratified = cross_val_score(
61     pipeline, X, y,
62     cv=skf,
63     scoring='roc_auc',
64     n_jobs=-1
65 )
66
67 print("\n=== Stratified K-Fold ===")
68 print(f"Fold AUC scores : {np.round(scores_stratified, 4)}")
69 print(f"Mean AUC          : {scores_stratified.mean():.4f}")
70 print(f"Std AUC          : {scores_stratified.std():.4f}")
71
72 # -----
73 # 5. cross_validate returns richer information:
74 #   train scores, test scores, and fit/score times.
75 # -----
76 cv_results = cross_validate(
77     pipeline, X, y,
78     cv=skf,
79     scoring=['roc_auc', 'f1', 'precision', 'recall'],
80     return_train_score=True, # detect overfitting by comparing
81     n_jobs=-1
82 )
83
84 results_df = pd.DataFrame({
85     'train_roc_auc': cv_results['train_roc_auc'],
86     'test_roc_auc' : cv_results['test_roc_auc'],
87     'test_f1'      : cv_results['test_f1'],
88     'fit_time_s'   : cv_results['fit_time']
89 })
90 print("\n=== Detailed CV Results ===")
91 print(results_df.round(4).to_string(index=False))
92
93 # If train AUC >> test AUC, the model is overfitting.
94 gap = cv_results['train_roc_auc'].mean() - cv_results['test_roc_auc']
95     .mean()
96 print(f"\nOverfit gap (train - test AUC): {gap:.4f}")
97 if gap > 0.05:
98     print("WARNING: Possible overfitting detected.")

```

What this code does, step by step:

Step 1 creates a synthetic binary classification dataset with mild class imbalance, which mirrors real business scenarios where one outcome (e.g., fraud, churn) is rarer than the other.

Step 2 is the most important architectural decision: the `Pipeline` ensures that the `StandardScaler` is fitted *only* on the training portion of each fold. If you scaled all data

first and then cross-validated, the validation fold's statistics would have influenced the scaler — a subtle form of data leakage that inflates your score.

Steps 3 and 4 compare plain `KFold` with `StratifiedKFold`. In practice, stratification reduces the variance of the CV estimator when class imbalance is present.

Step 5 uses

1.3 Time-Series Cross-Validation

Standard K-Fold **must not be used with time-series data**. The reason is fundamental: the future cannot cause the past. If your validation fold contains data from January and your training fold contains data from March, your model learns from the future to predict the past. This produces wildly optimistic estimates that will not hold in live deployment.

1.3.1 Walk-Forward Validation (Expanding Window)

`TimeSeriesSplit` implements *walk-forward validation*, also called *expanding window* or *rolling origin* cross-validation. The training window always ends before the validation window begins, and the training set grows with each split.

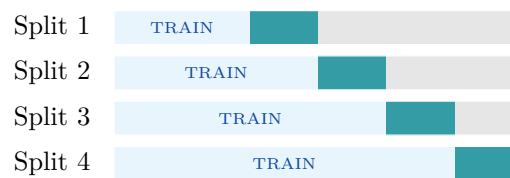


Figure 2: Time-Series Walk-Forward Validation: training always precedes validation temporally.

Time-Series Cross-Validation

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import TimeSeriesSplit, cross_validate
4 from sklearn.ensemble import GradientBoostingClassifier
5 from sklearn.pipeline import Pipeline
6 from sklearn.preprocessing import StandardScaler
7
8 # -----
9 # 1. Simulate a time-ordered dataset (e.g., daily transaction data)
10 #    n_samples = 2000 days; features represent behavioural signals.
11 # -----
12 np.random.seed(42)
13 n = 2000
14
15 # Create a temporal trend in the signal to mimic real time-series
16 time_idx = np.arange(n)
17 signal = 0.3 * np.sin(2 * np.pi * time_idx / 365) # seasonal
18           component
19 X = np.column_stack([
20     signal + np.random.randn(n) * 0.5, # feature 1: noisy signal
21     np.random.randn(n),                # feature 2: noise
22 ])
```



```

22     time_idx / n + np.random.randn(n)*0.1 # feature 3: trend
23 ] )
24 y = (signal + np.random.randn(n) * 0.4 > 0).astype(int)
25
26 # -----
27 # 2. Configure TimeSeriesSplit
28 #     n_splits=5 creates 5 folds, each with progressively more
29 #     training
30 #     data. gap=30 inserts a 30-day gap between training and
31 #     validation
32 #     to simulate realistic deployment delay.
33 # -----
34 tscv = TimeSeriesSplit(n_splits=5, gap=30)
35
36 # -----
37 # 3. Inspect what each split looks like
38 # -----
39 print("Split summary (train size | val size):")
40 for fold_num, (train_idx, val_idx) in enumerate(tscv.split(X), start
41 =1):
42     print(f"    Fold {fold_num}: train [{train_idx[0]}..{train_idx[-1]}] "
43           f"({len(train_idx)} samples) | "
44           f"val [{val_idx[0]}..{val_idx[-1]}] ({len(val_idx)} "
45           f"samples)")
46
47 # -----
48 # 4. Cross-validate a gradient boosting classifier
49 # -----
50 pipeline = Pipeline([
51     ('scaler', StandardScaler()),
52     ('clf', GradientBoostingClassifier(
53         n_estimators=100, learning_rate=0.05,
54         max_depth=4, random_state=42
55     ))
56 ])
57
58 cv_results = cross_validate(
59     pipeline, X, y,
60     cv=tscv,
61     scoring=['roc_auc', 'f1'],
62     return_train_score=True
63 )
64
65 print("\n=== Time-Series CV Results ===")
66 for fold_i in range(5):
67     print(f"    Fold {fold_i+1}: "
68           f"train AUC={cv_results['train_roc_auc'][fold_i]:.4f} "
69           f"val AUC={cv_results['test_roc_auc'][fold_i]:.4f} "
70           f"val F1={cv_results['test_f1'][fold_i]:.4f}")
71
72 print(f"\nMean val AUC : {cv_results['test_roc_auc'].mean():.4f} "
73       f"+/- {cv_results['test_roc_auc'].std():.4f}")
74
75 # -----
76 # 5. Trend analysis: does performance degrade over time?
77 #     In live systems, concept drift causes later folds to score
78 #     lower.

```

```
74 # -----  
75 fold_aucs = cv_results['test_roc_auc']  
76 trend = np.polyfit(range(len(fold_aucs)), fold_aucs, deg=1)[0]  
77 print(f"\nPerformance trend across folds: {trend:+.4f} per fold")  
78 if trend < -0.01:  
79     print("ALERT: Declining performance potential concept drift.")  
80 else:  
81     print("Performance stable across time folds.")
```

The **gap parameter** is critical for production systems. If your model is retrained monthly and predictions are served for the following 30 days, inserting a 30-day gap into your CV scheme ensures that your offline evaluation mimics live conditions. Without it, validation samples immediately adjacent to the training cutoff may share patterns (such as the same weekly seasonality cycle) that won't exist at actual deployment time.

Trend analysis across folds is a uniquely time-series concern. If AUC falls monotonically from fold 1 to fold 5, your model is experiencing *concept drift* — the statistical relationship between features and target is changing over time. This is a signal to explore re-training schedules or online learning approaches.

- **Shuffling data** before splitting — destroys temporal ordering and guarantees leakage.
- **Using aggregated features** (e.g., 30-day rolling mean computed on the full dataset) before splitting. Always compute such features inside the pipeline or after splitting.
- **Ignoring the deployment gap** between training cutoff and live prediction.

2 Hyperparameter Tuning: Systematic Model Configuration

2.1 What Are Hyperparameters?

Machine learning models have two kinds of parameters. *Model parameters* are learned during training (e.g., the weights in a neural network, the split thresholds in a decision tree). *Hyperparameters* are set *before* training begins and control the learning process itself: how many trees to grow, how deep each tree can be, the learning rate, regularisation strength, and so on.

The goal of hyperparameter tuning is to find the configuration that maximises the model's generalisation performance, as measured by cross-validation, not by training error.

Every tuning method follows the same loop: **propose** a configuration → **train** the model with that configuration → **evaluate** using cross-validation → **record** the score → repeat until budget exhausted → **select** the configuration with the best CV score. The methods differ in how they *propose* the next candidate.

2.2 Grid Search

Grid search exhaustively evaluates every combination of a pre-specified set of hyperparameter values. If you specify 3 values for `max_depth` and 4 values for `n_estimators`, grid search trains $3 \times 4 = 12$ models (each cross-validated k times), giving $12k$ total fits.

This is guaranteed to find the best configuration *within the grid*, but it scales exponentially with the number of hyperparameters – a phenomenon called the *curse of dimensionality in search*.

Grid Search with Cross-Validation

```

1  import numpy as np
2  import pandas as pd
3  from sklearn.datasets import make_classification
4  from sklearn.ensemble import GradientBoostingClassifier
5  from sklearn.model_selection import GridSearchCV, StratifiedKFold
6  from sklearn.pipeline import Pipeline
7  from sklearn.preprocessing import StandardScaler
8
9  X, y = make_classification(
10     n_samples=1500, n_features=20, n_informative=12,
11     weights=[0.65, 0.35], random_state=42
12 )
13
14  # -----
15  # Pipeline: scaler + classifier
16  # Hyperparameters for pipeline steps are referenced as
17  # "stepname__parametername"
18  # -----
19  pipeline = Pipeline([

```

```

20     ('scaler', StandardScaler()),
21     ('clf', GradientBoostingClassifier(random_state=42))
22 ])
23
24 # -----
25 # Define the grid
26 # 3 x 4 x 3 = 36 configurations, each with 5-fold CV = 180 fits.
27 # For production models, this grid is modest; real searches often
28 # involve hundreds of configurations.
29 # -----
30 param_grid = {
31     'clf__n_estimators': [50, 100, 200],
32     'clf__max_depth': [2, 3, 4, 5],
33     'clf__learning_rate': [0.01, 0.05, 0.1],
34 }
35
36 cv_strategy = StratifiedKFold(n_splits=5, shuffle=True, random_state
    =42)
37
38 grid_search = GridSearchCV(
39     estimator=pipeline,
40     param_grid=param_grid,
41     scoring='roc_auc',           # optimise for AUC
42     cv=cv_strategy,
43     n_jobs=-1,                  # parallelise across all cores
44     verbose=1,                  # print progress
45     refit=True                  # refit on full data with best params
46 )
47
48 grid_search.fit(X, y)
49
50 # -----
51 # Results exploration
52 # -----
53 print(f"Best CV AUC : {grid_search.best_score_:.4f}")
54 print(f"Best params : {grid_search.best_params_}")
55
56 # Convert the full results to a DataFrame for analysis
57 results_df = pd.DataFrame(grid_search.cv_results_)
58
59 # Show top-5 configurations
60 top5 = results_df[['params', 'mean_test_score', 'std_test_score',
61                    'rank_test_score']]\
62     .sort_values('rank_test_score').head(5)
63
64 print("\nTop 5 configurations:")
65 for _, row in top5.iterrows():
66     print(f"  Rank {int(row['rank_test_score'])}: "
67           f"AUC={row['mean_test_score']:.4f} "
68           f"+/-{row['std_test_score']:.4f} "
69           f"{row['params']}")
70
71 # -----
72 # The best estimator is already refitted on the full dataset.
73 # It can be used directly for prediction.
74 # -----
75 best_model = grid_search.best_estimator_

```

2.3 Random Search

Random search samples hyperparameter combinations uniformly at random from a defined distribution. Bergstra and Bengio (2012) demonstrated mathematically that random search is more efficient than grid search when the number of hyperparameters is large and only a few are truly important, because random search explores more of the *marginal* space of each parameter with the same number of evaluations.

The intuition: if 2 out of 5 hyperparameters matter but you don't know which 2, a $5 \times 5 \times 5 \times 5 \times 5$ grid evaluates each important parameter at only 5 distinct values. Fifty random samples explore each important parameter at up to 50 distinct values, while still covering the full joint space.

Random Search with Continuous Distributions

```

1  import numpy as np
2  from scipy.stats import uniform, randint, loguniform
3  from sklearn.ensemble import GradientBoostingClassifier
4  from sklearn.model_selection import RandomizedSearchCV,
   StratifiedKFold
5  from sklearn.pipeline import Pipeline
6  from sklearn.preprocessing import StandardScaler
7  from sklearn.datasets import make_classification
8
9  X, y = make_classification(
10     n_samples=1500, n_features=20, n_informative=12,
11     weights=[0.65, 0.35], random_state=42
12 )
13
14 pipeline = Pipeline([
15     ('scaler', StandardScaler()),
16     ('clf', GradientBoostingClassifier(random_state=42))
17 ])
18
19 # -----
20 # Distributions instead of grids
21 #
22 # loguniform(a, b): draws from a log-uniform distribution between a
23 #   and b, which is appropriate for learning rates and
24 #   regularisation
25 #   strengths that span orders of magnitude.
26 #
27 # randint(low, high): draws integers from [low, high).
28 #   Best for counts like n_estimators, max_depth.
29 #
30 # uniform(loc, scale): draws from [loc, loc+scale].
31 #   Best for continuous values with no strong scale preference.
32 # -----
33 param_distributions = {
34     'clf__n_estimators': randint(50, 500),          # integers in [
35     'clf__max_depth':   randint(2, 8),              # integers in [
36     'clf__learning_rate': loguniform(1e-3, 0.5),    # log-scale [0.
37     'clf__subsample':   uniform(0.5, 0.5),          # uniform in [0
38     '.5, 1.0]

```

```

37     'clf__min_samples_leaf': randint(1, 30),
38     'clf__max_features':    uniform(0.3, 0.7),           # fraction of
                        features
39 }
40
41 cv_strategy = StratifiedKFold(n_splits=5, shuffle=True, random_state
                        =42)
42
43 random_search = RandomizedSearchCV(
44     estimator=pipeline,
45     param_distributions=param_distributions,
46     n_iter=60,           # evaluate 60 random configurations
47     scoring='roc_auc',
48     cv=cv_strategy,
49     n_jobs=-1,
50     random_state=42,
51     verbose=1,
52     refit=True
53 )
54
55 random_search.fit(X, y)
56
57 print(f"Best CV AUC   : {random_search.best_score_:.4f}")
58 print(f"Best params   :")
59 for k, v in random_search.best_params_.items():
60     print(f"    {k:<35}: {v}")
61
62 # -----
63 # Compare grid search vs random search coverage
64 # For 60 evaluations:
65 #   Grid: constrained to pre-specified values
66 #   Random: explores continuous ranges, likely finds better optima
67 # -----

```

2.4 Bayesian Optimisation

Bayesian optimisation is the state of the art for expensive hyperparameter tuning. Unlike grid and random search, which treat each evaluation independently, Bayesian optimisation *learns* from previous evaluations and uses a probabilistic *surrogate model* (typically a Gaussian Process or Tree-structured Parzen Estimator) to predict which regions of the hyperparameter space are most promising.

The surrogate model balances *exploitation* (evaluating where the model predicts the best performance) with *exploration* (evaluating uncertain regions that might harbour better optima). This trade-off is managed by an *acquisition function*, most commonly Expected Improvement (EI):

$$\text{EI}(\mathbf{x}) = \mathbb{E}[\max(f(\mathbf{x}) - f(\mathbf{x}^+), 0)]$$

where $f(\mathbf{x}^+)$ is the best observed score so far. The next evaluation point is $\mathbf{x}_{\text{next}} = \arg \max_{\mathbf{x}} \text{EI}(\mathbf{x})$.

Bayesian Optimisation with scikit-optimize

```

1  import numpy as np
2  from sklearn.datasets import make_classification
3  from sklearn.ensemble import GradientBoostingClassifier
4  from sklearn.model_selection import StratifiedKFold, cross_val_score
5  from sklearn.pipeline import Pipeline
6  from sklearn.preprocessing import StandardScaler
7
8  # scikit-optimize provides BayesSearchCV with a scikit-learn API
9  # Install: pip install scikit-optimize
10 from skopt import BayesSearchCV
11 from skopt.space import Real, Integer, Categorical
12
13 X, y = make_classification(
14     n_samples=1500, n_features=20, n_informative=12,
15     weights=[0.65, 0.35], random_state=42
16 )
17
18 pipeline = Pipeline([
19     ('scaler', StandardScaler()),
20     ('clf', GradientBoostingClassifier(random_state=42))
21 ])
22
23 # -----
24 # Search space using skopt types
25 # Real(low, high, prior='log-uniform'): continuous log-scale range
26 # Integer(low, high): integer range
27 # Categorical([...]): discrete set of choices
28 # -----
29 search_space = {
30     'clf__n_estimators' : Integer(50, 400),
31     'clf__max_depth'    : Integer(2, 7),
32     'clf__learning_rate': Real(1e-3, 0.5, prior='log-uniform'),
33     'clf__subsample'    : Real(0.4, 1.0),
34     'clf__min_samples_leaf': Integer(1, 25),
35 }
36
37 cv_strategy = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
38
39 # -----
40 # BayesSearchCV: n_iter controls total evaluations.
41 # Fewer evaluations than grid/random search can achieve better
42 #   results
43 # because each evaluation is informed by prior results.
44 # -----
45 bayes_search = BayesSearchCV(
46     estimator=pipeline,
47     search_spaces=search_space,
48     n_iter=40,           # 40 informed evaluations
49     scoring='roc_auc',
50     cv=cv_strategy,
51     n_jobs=1,           # skopt internals require n_jobs=1
52     random_state=42,
53     verbose=0,
54     refit=True
55 )

```

```

55
56 bayes_search.fit(X, y)
57
58 print(f"Best Bayesian CV AUC : {bayes_search.best_score_:.4f}")
59 print("Best parameters:")
60 for k, v in bayes_search.best_params_.items():
61     print(f"    {k:<35}: {v}")
62
63 # -----
64 # Convergence plot: track how AUC improves iteration by iteration.
65 # A flat curve means the search has converged; more budget won't
66 #   help.
67 # -----
68
69 import matplotlib.pyplot as plt
70
71 scores_over_time = []
72 best_so_far = -np.inf
73 for score in bayes_search.cv_results_['mean_test_score']:
74     best_so_far = max(best_so_far, score)
75     scores_over_time.append(best_so_far)
76
77 plt.figure(figsize=(8, 4))
78 plt.plot(range(1, len(scores_over_time)+1), scores_over_time,
79         marker='o', color='#0d47a1', linewidth=2, markersize=4)
80 plt.xlabel("Evaluation number")
81 plt.ylabel("Best CV AUC so far")
82 plt.title("Bayesian Optimisation Convergence")
83 plt.grid(alpha=0.3)
84 plt.tight_layout()
85 plt.savefig("bayesian_convergence.png", dpi=150)
86 plt.close()
87 print("Convergence plot saved.")

```

2.4.1 Choosing a Tuning Strategy

The following table summarises when to use each approach:

Method	Best for	Weakness	Budget
Grid Search	Small grids, interpretable search space	Exponential scaling	Low
Random Search	Many hyperparameters, few matter	May miss fine structure	Medium
Bayesian Opt.	Expensive models, large spaces	Sequential (harder to parallelise)	High

Table 1: Hyperparameter tuning method comparison.

When you tune hyperparameters with 5-fold stratified CV, you must also report your final performance estimate using a separate *outer* cross-validation loop (nested CV), or a hold-out test set that was never touched during tuning. Using the same data for both tuning and final reporting is a form of overfitting.

3 Metrics for Imbalanced Datasets

3.1 Why Accuracy Fails

Imagine a fraud detection model. Fraudulent transactions represent 0.5% of all transactions. A model that predicts “not fraud” for every single transaction achieves 99.5% accuracy. This is a useless model that catches zero fraud.

This is why *class imbalance* is one of the most practically important challenges in applied machine learning. The choice of evaluation metric must reflect the cost structure of the problem, not just raw prediction accuracy.

3.2 The Confusion Matrix: A Foundation

All classification metrics derive from four counts:

	Predicted Positive	Predicted Negative
Actual Positive	green!15 TP (True Positive)	red!15 FN (False Negative)
Actual Negative	red!15 FP (False Positive)	green!15 TN (True Negative)

Table 2: The 2CE2 confusion matrix.

From these four numbers, we derive every metric in this section.

3.3 Core Metrics

Precision answers the question: “Of all the cases I predicted positive, what fraction actually were positive?” It penalises false alarms.

$$\text{Precision} = \frac{TP}{TP + FP}$$

Recall (also called Sensitivity or True Positive Rate) answers: “Of all the actual positive cases, what fraction did I catch?” It penalises missed detections.

$$\text{Recall} = \frac{TP}{TP + FN}$$

F1-Score is the harmonic mean of precision and recall. The harmonic mean is used (rather than arithmetic) because it punishes extreme imbalances between precision and recall.

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

The F_β score generalises F1 by weighting recall β times more than precision. For fraud detection where missing a fraud (false negative) is far more costly than a false alarm, use $\beta > 1$:

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}$$

Balanced Accuracy averages recall across all classes, correcting for imbalance without requiring threshold tuning:

$$\text{Balanced Accuracy} = \frac{1}{2} \left(\frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right)$$

Comprehensive Imbalanced Classification Evaluation

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  from sklearn.datasets import make_classification
4  from sklearn.ensemble import GradientBoostingClassifier,
   RandomForestClassifier
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.model_selection import StratifiedKFold, cross_val_
   predict
7  from sklearn.pipeline import Pipeline
8  from sklearn.preprocessing import StandardScaler
9  from sklearn.metrics import (
10     classification_report,
11     confusion_matrix,
12     f1_score,
13     balanced_accuracy_score,
14     precision_recall_curve,
15     average_precision_score,
16     roc_auc_score,
17     ConfusionMatrixDisplay
18 )
19
20 # -----
21 # 1. Heavily imbalanced dataset: 95% class 0, 5% class 1
22 #    This simulates fraud, disease, or rare equipment failure.
23 # -----
24 X, y = make_classification(
25     n_samples=5000,
26     n_features=20,
27     n_informative=10,
28     weights=[0.95, 0.05], # 4750 negatives, 250 positives
29     random_state=42
30 )
31 print(f"Class distribution: {np.bincount(y)}  "
32       f"({100*y.mean():.1f}% positive)")
33
34 # -----
35 # 2. Train three models and get out-of-fold predictions
36 #    cross_val_predict generates predictions for every sample using
37 #    only models trained without that sample proper OOF evaluation.
38 # -----

```

```

39 models = {
40     'Logistic Regression': Pipeline([
41         ('scaler', StandardScaler()),
42         ('clf', LogisticRegression(
43             class_weight='balanced', # upweight minority
44             C=0.1, max_iter=1000
45         ))
46     ]),
47     'Random Forest': Pipeline([
48         ('scaler', StandardScaler()),
49         ('clf', RandomForestClassifier(
50             n_estimators=200,
51             class_weight='balanced_subsample',
52             random_state=42
53         ))
54     ]),
55     'Gradient Boosting': Pipeline([
56         ('scaler', StandardScaler()),
57         ('clf', GradientBoostingClassifier(
58             n_estimators=200, learning_rate=0.05,
59             max_depth=4, random_state=42
60         ))
61     ]),
62 }
63
64 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
65
66 results = {}
67 for name, pipe in models.items():
68     # predict_proba: probability scores for each class
69     proba = cross_val_predict(pipe, X, y, cv=cv, method='predict_
70                               proba')[0,1]
71     # predict: hard class labels at 0.5 threshold
72     labels = cross_val_predict(pipe, X, y, cv=cv, method='predict')
73
74     results[name] = {
75         'proba' : proba,
76         'labels' : labels,
77         'ap' : average_precision_score(y, proba), # area under
78             P-R curve
79         'auc' : roc_auc_score(y, proba),
80         'f1' : f1_score(y, labels),
81         'f2' : f1_score(y, labels, beta=2.0), # recall-
82             weighted
83         'bal_acc': balanced_accuracy_score(y, labels),
84     }
85
86 # -----
87 # 3. Summary table
88 # -----
89 print("\n{<22} {>8} {>8} {>8} {>8} {>8}".format(
90     "Model", "AUC", "Avg Prec", "F1", "F2", "BalAcc"))
91 print("-" * 70)
92 for name, r in results.items():
93     print(f"{name:<22} {r['auc']:>8.4f} {r['ap']:>8.4f} "
94           f"{r['f1']:>8.4f} {r['f2']:>8.4f} {r['bal_acc']:>8.4f}")
95
96 # -----

```

```

94 # 4. Precision-Recall curves the right plot for imbalanced data
95 # ROC curves are optimistic when the negative class dominates;
96 # P-R curves expose the real difficulty of finding true positives
97 # -----
98 fig, axes = plt.subplots(1, 2, figsize=(13, 5))
99
100 colors = ['#0d47a1', '#00838f', '#e65100']
101 for (name, r), color in zip(results.items(), colors):
102     prec, rec, thresh = precision_recall_curve(y, r['proba'])
103     axes[0].plot(rec, prec, color=color, linewidth=2,
104                 label=f"{name} (AP={r['ap']:.3f})")
105
106 axes[0].axhline(y.mean(), color='gray', linestyle='--', linewidth=1.
107                2,
108                label=f"Random (AP={y.mean():.3f})")
109 axes[0].set_xlabel("Recall", fontsize=12)
110 axes[0].set_ylabel("Precision", fontsize=12)
111 axes[0].set_title("Precision-Recall Curves", fontsize=13, fontweight
112                  = 'bold')
113 axes[0].legend(fontsize=9)
114 axes[0].grid(alpha=0.3)
115
116 # -----
117 # 5. Threshold analysis for the best model
118 # The default 0.5 threshold is rarely optimal for imbalanced data
119 # Plot F1 and F2 as a function of threshold to find the sweet
120 # spot.
121 # -----
122 best_name = max(results, key=lambda n: results[n]['ap'])
123 proba_best = results[best_name]['proba']
124 prec, rec, thresholds = precision_recall_curve(y, proba_best)
125
126 f1_scores = 2 * prec[:-1] * rec[:-1] / (prec[:-1] + rec[:-1] + 1e-9)
127 f2_scores = 5 * prec[:-1] * rec[:-1] / (4*prec[:-1] + rec[:-1] + 1e-
128    9)
129
130 axes[1].plot(thresholds, f1_scores, color='#0d47a1', label='F1 score
131    ', lw=2)
132 axes[1].plot(thresholds, f2_scores, color='#e65100', label='F2 score
133    ', lw=2)
134 axes[1].axvline(0.5, color='gray', linestyle='--', label='Default
135    threshold')
136
137 best_f1_thresh = thresholds[np.argmax(f1_scores)]
138 axes[1].axvline(best_f1_thresh, color='#0d47a1', linestyle=':', lw=1
139    .5,
140    label=f'Best F1 threshold={best_f1_thresh:.2f}')
141
142 axes[1].set_xlabel("Decision Threshold", fontsize=12)
143 axes[1].set_ylabel("Score", fontsize=12)
144 axes[1].set_title(f"F1/F2 vs Threshold {best_name}", fontsize=13,
145                  fontweight='bold')
146 axes[1].legend(fontsize=9)
147 axes[1].grid(alpha=0.3)
148 axes[1].set_xlim(0, 1)
149

```

```

142 plt.tight_layout()
143 plt.savefig("imbalanced_metrics.png", dpi=150)
144 plt.close()
145
146 # -----
147 # 6. Optimal threshold classification report
148 # -----
149 optimal_labels = (proba_best >= best_f1_thresh).astype(int)
150 print(f"\nClassification report at optimal threshold ({best_f1_
      thresh:.2f}):")
151 print(classification_report(y, optimal_labels,
152                             target_names=['No Fraud', 'Fraud']))

```

3.3.1 Reading Precision-Recall Curves

The *Precision-Recall (PR) curve* plots precision on the y-axis against recall on the x-axis as the decision threshold sweeps from 0 to 1. The *Average Precision (AP)* score is the area under this curve.

For imbalanced datasets, the PR curve is more informative than the ROC curve because it directly measures performance on the minority class. The ROC curve can appear excellent ($AUC > 0.9$) even when the model performs poorly at actually finding positive cases, because the large number of true negatives inflates the TN count.

The dashed horizontal line at the baseline class frequency is the AP score of a random classifier. Any model must substantially exceed this baseline to be useful.

3.3.2 Threshold Optimisation

The decision threshold of 0.5 is an arbitrary convention. For imbalanced problems, the optimal threshold is almost never 0.5. The threshold should be chosen based on the cost asymmetry of the problem: if false negatives are 5 times more costly than false positives, use $\beta = \sqrt{5} \approx 2.2$ in the F_β score to find the optimal threshold.

Metric	Use when	Range
Accuracy	Classes are balanced AND errors are equally costly	[0, 1]
Balanced Accuracy	Imbalanced, quick threshold-free comparison	[0, 1]
F1 Score	False positives and negatives equally costly	[0, 1]
F2 Score	Missing positives more costly than false alarms	[0, 1]
Average Precision	Need full threshold-independent evaluation	[0, 1]
ROC-AUC	Comparing discriminative ability across thresholds	[0.5, 1]

Table 3: Metric selection guide for imbalanced classification problems.

4 Model Interpretability

4.1 Why Interpretability Is Not Optional

A model that cannot be explained is a model that cannot be trusted, audited, debugged, or improved. In a business context, interpretability matters for three distinct reasons:

Regulatory compliance: Legislation such as GDPR (right to explanation), the EU AI Act, and US financial regulations often require that models provide explanations for decisions that affect individuals.

Model debugging: If your model performs worse on a demographic subgroup or in a specific time period, interpretability methods reveal which features are driving that behaviour.

Stakeholder trust: A data scientist who can say “the model flags this customer as high-risk primarily because their last payment was 60 days late and their credit utilisation is 94%” earns far more trust than one who shows a confusion matrix.

4.2 SHAP: SHapley Additive exPlanations

4.2.1 The Theory

SHAP is grounded in cooperative game theory. A *Shapley value* is the average marginal contribution of a feature across all possible orderings (coalitions) of features. For feature j in a prediction $f(\mathbf{x})$:

$$\phi_j = \sum_{S \subseteq F \setminus \{j\}} \frac{|S|! (|F| - |S| - 1)!}{|F|!} [f_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) - f_S(\mathbf{x}_S)]$$

where F is the set of all features, S is any subset not containing j , and f_S is the model’s prediction using only features in S . The SHAP value ϕ_j represents the feature’s fair share of the difference between the model’s prediction and the global base rate.

The key property that makes SHAP uniquely appealing is that the SHAP values sum exactly to the model’s output relative to the expected output:

$$f(\mathbf{x}) = \mathbb{E}[f(X)] + \sum_{j=1}^{|F|} \phi_j$$

This additive decomposition means every prediction is fully accounted for, with no residual “unexplained” component.

SHAP: Global and Local Explanations

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib
4 matplotlib.use('Agg')
5 import matplotlib.pyplot as plt
6 import shap
```

```

7  from sklearn.datasets import make_classification
8  from sklearn.ensemble import GradientBoostingClassifier
9  from sklearn.model_selection import train_test_split
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.pipeline import Pipeline
12
13 # -----
14 # 1. Dataset with interpretable feature names
15 #   (simulating a credit scoring context)
16 # -----
17 np.random.seed(42)
18 n = 2000
19 X_raw = pd.DataFrame({
20     'days_late_last_payment' : np.random.exponential(10, n).clip(0,
21         120),
22     'credit_utilisation_pct' : np.random.beta(2, 5, n) * 100,
23     'num_open_accounts'      : np.random.poisson(4, n),
24     'income_annual_k'        : np.random.lognormal(4, 0.5, n),
25     'years_credit_history'    : np.random.gamma(3, 3, n).clip(0.5,
26         30),
27     'num_hard_inquiries'      : np.random.poisson(1.5, n),
28     'debt_to_income_ratio'    : np.random.beta(2, 4, n),
29     'months_since_derog'      : np.random.exponential(20, n).clip(0,
30         120),
31 })
32
33 # Logistic rule with noise creates realistic credit default target
34 logit = (
35     0.05 * X_raw['days_late_last_payment']
36     + 0.03 * X_raw['credit_utilisation_pct']
37     - 0.02 * X_raw['income_annual_k'] / 10
38     + 0.5 * X_raw['debt_to_income_ratio']
39     + 0.4 * np.random.randn(n)
40     - 3.0
41 )
42 y = (1 / (1 + np.exp(-logit)) > 0.5).astype(int)
43
44 X_train, X_test, y_train, y_test = train_test_split(
45     X_raw, y, test_size=0.25, random_state=42, stratify=y
46 )
47
48 # -----
49 # 2. Train a gradient boosting model
50 #   We do NOT use a Pipeline here because SHAP needs the raw
51 #   model object (not wrapped), though TreeExplainer handles this.
52 # -----
53 scaler = StandardScaler()
54 X_train_s = scaler.fit_transform(X_train)
55 X_test_s = scaler.transform(X_test)
56
57 model = GradientBoostingClassifier(
58     n_estimators=200, learning_rate=0.05,
59     max_depth=4, random_state=42
60 )
61 model.fit(X_train_s, y_train)
62
63 # -----
64 # 3. Create SHAP explainer

```



```

62 # TreeExplainer is optimised for tree-based models (O(TLD^2)
    # where
63 # T=trees, L=leaves, D=depth) much faster than brute-force.
64 # -----
65 explainer = shap.TreeExplainer(model)
66
67 # Compute SHAP values for the test set
68 # shap_values shape: (n_test_samples, n_features) for binary
69 # For binary, index [1] gives SHAP for the positive class.
70 shap_values = explainer.shap_values(X_test_s)
71
72 # -----
73 # 4. Global explanation: Beeswarm summary plot
74 # Each dot = one test sample. X-axis = SHAP value (impact on
75 # model output). Colour = feature value (red=high, blue=low).
76 # Features are sorted by mean absolute SHAP value.
77 # -----
78 shap.summary_plot(
79     shap_values, X_test, # use un-scaled X for readable values
80     plot_type='dot',
81     show=False,
82     max_display=8
83 )
84 plt.tight_layout()
85 plt.savefig("shap_beeswarm.png", dpi=150, bbox_inches='tight')
86 plt.close()
87
88 # -----
89 # 5. Global feature importance (bar chart version)
90 # -----
91 shap.summary_plot(
92     shap_values, X_test,
93     plot_type='bar',
94     show=False,
95     max_display=8
96 )
97 plt.tight_layout()
98 plt.savefig("shap_importance.png", dpi=150, bbox_inches='tight')
99 plt.close()
100
101 # -----
102 # 6. Local explanation: a single prediction
103 # Interpret the decision for one high-risk customer.
104 # -----
105 sample_idx = np.where(y_test == 1)[0][0] # first true positive
106 sample = X_test_s[sample_idx:sample_idx+1]
107
108 predicted_prob = model.predict_proba(sample)[0, 1]
109 print(f"\n--- Local Explanation for Sample {sample_idx} ---")
110 print(f"Predicted probability of default: {predicted_prob:.3f}")
111 print(f"Base value (mean prediction) : {explainer.expected_value:
    .3f}")
112 print(f"\nSHAP contribution of each feature:")
113
114 feature_names = X_raw.columns.tolist()
115 sample_shap = shap_values[sample_idx]
116 sample_orig = X_test.iloc[sample_idx]
117

```

```

118 contributions = pd.DataFrame({
119     'feature'      : feature_names,
120     'value'        : sample_orig.values,
121     'shap_value'   : sample_shap
122 }).sort_values('shap_value', key=abs, ascending=False)
123
124 for _, row in contributions.iterrows():
125     direction = "increases" if row['shap_value'] > 0 else "decreases"
126     print(f" {row['feature']:<30}: {direction} risk by "
127           f"{abs(row['shap_value']):.4f} "
128           f"(feature value = {row['value']:.2f})")
129
130 total = explainer.expected_value + sample_shap.sum()
131 print(f"\nBase rate ({explainer.expected_value:.3f}) + "
132       f"sum of contributions ({sample_shap.sum():.3f}) "
133       f"= {total:.3f} (model output before sigmoid)")
134
135 # -----
136 # 7. SHAP Dependence Plot: non-linear relationships
137 # Shows SHAP value of one feature as a function of its raw value,
138 # coloured by an interacting feature chosen automatically by SHAP
139 # -----
140 shap.dependence_plot(
141     'days_late_last_payment',
142     shap_values, X_test,
143     interaction_index='auto',
144     show=False
145 )
146 plt.tight_layout()
147 plt.savefig("shap_dependence.png", dpi=150, bbox_inches='tight')
148 plt.close()
149 print("\nSHAP plots saved.")

```

4.2.2 Reading SHAP Output

The **beeswarm plot** is the primary diagnostic for global interpretability. Each row is a feature. Each dot is a sample. The position on the x-axis is the SHAP value: positive values push the prediction toward the positive class; negative values push it away. The colour represents the feature's raw value. For example: if points with high `credit_utilisation_pct` (red) have large positive SHAP values, the model has learned that high utilisation predicts default — a result that is consistent with domain knowledge, which builds confidence.

The **local waterfall** for a single prediction shows how each feature's SHAP value adds to or subtracts from the base rate to produce the final score. This is the explanation you would present to a credit officer reviewing a declined application.

4.3 LIME: Local Interpretable Model-agnostic Explanations

LIME takes a philosophically different approach. Rather than computing exact Shapley values, LIME constructs a *local linear approximation* of the model's behaviour in the neighbourhood of a specific prediction.

The algorithm works in three steps:

- Step 1. Perturb the input.** Generate N synthetic samples by randomly turning features on and off (for tabular data, this means replacing feature values with values sampled from the training distribution).
- Step 2. Weight by proximity.** Weight each synthetic sample by its distance from the original input using an exponential kernel.
- Step 3. Fit a sparse linear model.** Train a regularised linear model (typically LASSO) on the weighted synthetic samples. The coefficients of this linear model are the local explanation.

LIME Tabular Explainer

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib
4  matplotlib.use('Agg')
5  import matplotlib.pyplot as plt
6  import lime
7  import lime.lime_tabular
8  from sklearn.ensemble import GradientBoostingClassifier
9  from sklearn.model_selection import train_test_split
10 from sklearn.preprocessing import StandardScaler
11
12 # Reuse the credit scoring dataset from the SHAP section
13 np.random.seed(42)
14 n = 2000
15 X_raw = pd.DataFrame({
16     'days_late_last_payment' : np.random.exponential(10, n).clip(0,
17     120),
18     'credit_utilisation_pct' : np.random.beta(2, 5, n) * 100,
19     'num_open_accounts' : np.random.poisson(4, n),
20     'income_annual_k' : np.random.lognormal(4, 0.5, n),
21     'years_credit_history' : np.random.gamma(3, 3, n).clip(0.5,
22     30),
23     'num_hard_inquiries' : np.random.poisson(1.5, n),
24     'debt_to_income_ratio' : np.random.beta(2, 4, n),
25     'months_since_derog' : np.random.exponential(20, n).clip(0,
26     120),
27 })
28 logit = (0.05*X_raw['days_late_last_payment']
29     + 0.03*X_raw['credit_utilisation_pct']
30     - 0.02*X_raw['income_annual_k']/10
31     + 0.5*X_raw['debt_to_income_ratio']
32     + 0.4*np.random.randn(n) - 3.0)
33 y = (1/(1+np.exp(-logit)) > 0.5).astype(int)
34
35 X_train, X_test, y_train, y_test = train_test_split(
36     X_raw.values, y, test_size=0.25, random_state=42, stratify=y
37 )
38 scaler = StandardScaler()
39 X_tr_s = scaler.fit_transform(X_train)
40 X_te_s = scaler.transform(X_test)

```

```

40 model = GradientBoostingClassifier(
41     n_estimators=200, learning_rate=0.05, max_depth=4, random_state=
42     42
43 )
44 model.fit(X_tr_s, y_train)
45
46 # -----
47 # 1. Build the LIME explainer
48 #     training_data: the SCALED training set, so LIME samples from
49 #         the
50 #         correct distribution.
51 #     feature_names: column names for readable output.
52 #     class_names: label strings.
53 #     discretize_continuous: if True, LIME bins continuous features
54 #         this makes explanations more readable but less precise.
55 # -----
56 explainer = lime.lime_tabular.LimeTabularExplainer(
57     training_data=X_tr_s,
58     feature_names=X_raw.columns.tolist(),
59     class_names=['No Default', 'Default'],
60     mode='classification',
61     discretize_continuous=True,
62     random_state=42
63 )
64
65 # -----
66 # 2. Explain a single prediction
67 #     num_features: how many features to include in the local model.
68 #     num_samples: how many perturbed samples to generate (more =
69 #         stable).
70 # -----
71 idx = np.where(y_test == 1)[0][3]    # a true default case
72 sample = X_te_s[idx]
73
74 explanation = explainer.explain_instance(
75     data_row=sample,
76     predict_fn=model.predict_proba,
77     num_features=8,
78     num_samples=3000,
79     top_labels=1
80 )
81
82 pred_prob = model.predict_proba(sample.reshape(1,-1))[0,1]
83 print(f"Predicted probability of default: {pred_prob:.3f}")
84 print("\nLIME local explanation (feature: weight):")
85 for feat, weight in explanation.as_list(label=1):
86     direction = "increases risk" if weight > 0 else "decreases risk"
87     print(f"    {feat:<50}: {weight:+.4f}    ({direction})")
88
89 # -----
90 # 3. Save the LIME explanation as an HTML file
91 #     This is LIME's native format and renders feature weights as
92 #     a horizontal bar chart.
93 # -----
94 html_str = explanation.as_html()
95 with open("lime_explanation.html", "w") as f:
96     f.write(html_str)
97 print("\nLIME HTML explanation saved.")

```

```

95
96 # -----
97 # 4. Manual bar chart (for embedding in reports)
98 # -----
99 lime_list = sorted(explanation.as_list(label=1), key=lambda x: x[1])
100 features_l = [item[0] for item in lime_list]
101 weights_l = [item[1] for item in lime_list]
102 colors_l = ['#e65100' if w > 0 else '#0d47a1' for w in weights_l]
103
104 fig, ax = plt.subplots(figsize=(9, 5))
105 bars = ax.barh(features_l, weights_l, color=colors_l, edgecolor='
    white')
106 ax.axvline(0, color='black', linewidth=0.8)
107 ax.set_xlabel("LIME weight (positive = increases P(Default))",
    fontsize=11)
108 ax.set_title(f"LIME Local Explanation Sample {idx}\n"
109             f"P(Default)={pred_prob:.3f}", fontsize=12, fontweight=
    'bold')
110 ax.grid(axis='x', alpha=0.3)
111 plt.tight_layout()
112 plt.savefig("lime_explanation.png", dpi=150, bbox_inches='tight')
113 plt.close()
114 print("LIME bar chart saved.")

```

4.3.1 SHAP vs LIME: When to Use Each

Property	SHAP	LIME
Mathematical foundation	Shapley values (game theory)	Local linear surrogate
Consistency	Globally consistent; local explanations sum to prediction	Locally consistent only; can vary between runs
Model agnosticism	TreeExplainer requires tree models; KernelExplainer is model-agnostic but slow	Fully model-agnostic
Computational cost	Fast for trees, slow for deep models	Moderate; depends on num_samples
Stability	Deterministic (for TreeExplainer)	Stochastic; run multiple times to check variance
Best for	Audits, global analysis, regulatory documentation	Quick local explanations, non-tree models

Table 4: SHAP versus LIME comparison.

4.4 Partial Dependence Plots

Partial Dependence Plots (PDP) show the *marginal* effect of one or two features on the model's predicted outcome, averaging over all other features. Formally, for feature j :

$$\hat{f}_j(x_j) = \mathbb{E}_{X_{-j}} \left[\hat{f}(x_j, X_{-j}) \right] \approx \frac{1}{n} \sum_{i=1}^n \hat{f}(x_j, \mathbf{x}_{-j}^{(i)})$$

For each value of x_j , the model prediction is computed for every training sample (fixing x_j while using each sample's actual values for all other features), and the results are averaged. This removes the dependence on all other features, showing the isolated effect of x_j .

Partial Dependence Plots and ICE Curves

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib
4  matplotlib.use('Agg')
5  import matplotlib.pyplot as plt
6  from sklearn.inspection import PartialDependenceDisplay
7  from sklearn.ensemble import GradientBoostingClassifier
8  from sklearn.model_selection import train_test_split
9  from sklearn.preprocessing import StandardScaler
10
11 np.random.seed(42)
12 n = 2000
13 X_raw = pd.DataFrame({
14     'days_late_last_payment' : np.random.exponential(10, n).clip(0,
15         120),
16     'credit_utilisation_pct' : np.random.beta(2, 5, n) * 100,
17     'num_open_accounts' : np.random.poisson(4, n),
18     'income_annual_k' : np.random.lognormal(4, 0.5, n),
19     'years_credit_history' : np.random.gamma(3, 3, n).clip(0.5,
20         30),
21     'num_hard_inquiries' : np.random.poisson(1.5, n),
22     'debt_to_income_ratio' : np.random.beta(2, 4, n),
23     'months_since_derog' : np.random.exponential(20, n).clip(0,
24         120),
25 })
26 logit = (0.05*X_raw['days_late_last_payment']
27     + 0.03*X_raw['credit_utilisation_pct']
28     - 0.02*X_raw['income_annual_k']/10
29     + 0.5*X_raw['debt_to_income_ratio']
30     + 0.4*np.random.randn(n) - 3.0)
31 y = (1/(1+np.exp(-logit)) > 0.5).astype(int)
32
33 X_train, X_test, y_train, y_test = train_test_split(
34     X_raw, y, test_size=0.25, random_state=42, stratify=y
35 )
36
37 scaler = StandardScaler()
38 X_tr_s = pd.DataFrame(scaler.fit_transform(X_train), columns=X_raw
39     .columns)
40 X_te_s = pd.DataFrame(scaler.transform(X_test), columns=X_raw.
41     columns)
42
43 model = GradientBoostingClassifier(
44     n_estimators=200, learning_rate=0.05, max_depth=4, random_state=
45     42
46 )
47 model.fit(X_tr_s, y_train)

```

```

42
43 # -----
44 # 1. PDP for the top 4 features
45 #     kind='both': overlay ICE (individual conditional expectation)
46 #     curves on top of the mean PDP line.
47 #     ICE curves reveal heterogeneity: if all ICE curves are parallel
48 #     ,
49 #     the feature has the same effect on everyone. If they fan out,
50 #     there are interaction effects with other features.
51 # -----
52 features_to_plot = [
53     'days_late_last_payment',
54     'credit_utilisation_pct',
55     'debt_to_income_ratio',
56     'income_annual_k'
57 ]
58
59 fig, axes = plt.subplots(2, 2, figsize=(13, 9))
60 axes_flat = axes.flatten()
61
62 for ax, feat in zip(axes_flat, features_to_plot):
63     disp = PartialDependenceDisplay.from_estimator(
64         model,
65         X_tr_s,
66         features=[feat],
67         kind='both',                # PDP + ICE
68         subsample=200,             # use 200 samples for ICE (speed)
69         n_jobs=-1,
70         grid_resolution=60,
71         ice_lines_kw={'color': '#90caf9', 'alpha': 0.3, 'linewidth':
72                     0.5},
73         pd_line_kw={'color': '#0d47a1', 'linewidth': 2.5, 'label': '
74                     PDP mean'})
75     ax.set_title(f"PDP + ICE: {feat}", fontsize=11, fontweight='bold')
76     ax.set_xlabel(feat, fontsize=9)
77     ax.set_ylabel("Partial dependence", fontsize=9)
78     ax.grid(alpha=0.3)
79
80 plt.suptitle("Partial Dependence Plots with ICE Curves",
81             fontsize=14, fontweight='bold', y=1.01)
82 plt.tight_layout()
83 plt.savefig("pdp_ice.png", dpi=150, bbox_inches='tight')
84 plt.close()
85
86 # -----
87 # 2. 2D PDP: interaction between two features
88 #     If the surface is not a simple sum of two 1D functions, there
89 #     is a genuine statistical interaction between the two features.
90 # -----
91 fig, ax = plt.subplots(figsize=(8, 6))
92 disp_2d = PartialDependenceDisplay.from_estimator(
93     model,
94     X_tr_s,
95     features=[('days_late_last_payment', 'credit_utilisation_pct')],
96     kind='average',

```

```
96     grid_resolution=30,
97     ax=ax,
98     n_jobs=-1
99 )
100 ax.set_title("2D PDP: Late Payment vs Credit Utilisation",
101              fontsize=12, fontweight='bold')
102 plt.tight_layout()
103 plt.savefig("pdp_2d.png", dpi=150, bbox_inches='tight')
104 plt.close()
105 print("PDP plots saved.")
```

Individual Conditional Expectation (ICE) curves expose something that PDPs hide. A PDP averages over all samples; if some samples see a positive effect from a feature and others see a negative effect, those effects cancel out and the PDP appears flat. ICE curves plot the predicted outcome for each individual sample as the feature varies, revealing this heterogeneity. When ICE curves cross each other, a feature interaction is present.

The **2D PDP** extends this to pairs of features. The contour plot shows the model's predicted outcome as a function of two features simultaneously. Non-linear contours that cannot be decomposed into a sum of a row pattern and a column pattern indicate an interaction effect.

Interpretability methods are not only for explainability to stakeholders — they are essential debugging tools. If a SHAP summary plot shows a feature contributing positively that you know from domain expertise should contribute negatively, you have found a bug: possibly a data leakage issue, a target encoding done incorrectly, or a labelling error. Always validate model behaviour against domain knowledge before deployment.

5 Putting It All Together: End-to-End Evaluation Workflow

This section demonstrates a complete, production-grade evaluation workflow that combines all four pillars from a single coherent script. The pattern shown here is the template you should follow for any classification problem.

Full Evaluation Pipeline — Production Template

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib
4  matplotlib.use('Agg')
5  import matplotlib.pyplot as plt
6  import shap
7  import warnings
8  warnings.filterwarnings('ignore')
9
10 from sklearn.datasets import make_classification
11 from sklearn.ensemble import GradientBoostingClassifier
12 from sklearn.model_selection import (
13     StratifiedKFold, cross_validate, RandomizedSearchCV, cross_val_
14     predict
15 )
16 from sklearn.pipeline import Pipeline
17 from sklearn.preprocessing import StandardScaler
18 from sklearn.metrics import (
19     classification_report, average_precision_score,
20     roc_auc_score, balanced_accuracy_score, f1_score,
21     precision_recall_curve
22 )
23 from scipy.stats import randint, loguniform, uniform
24
25 print("=" * 65)
26 print("  MODEL EVALUATION & VALIDATION  PRODUCTION WORKFLOW")
27 print("=" * 65)
28
29 # STEP 1: Data
30 #
31 X, y = make_classification(
32     n_samples=3000, n_features=20, n_informative=10,
33     weights=[0.88, 0.12], random_state=42
34 )
35 feature_names = [f"feature_{i:02d}" for i in range(X.shape[1])]
36 print(f"\nStep 1 Dataset: {X.shape[0]} samples, "
37       f"{X.shape[1]} features, {100*y.mean():.1f}% positive")
38
39 # Reserve a true holdout set touch this ONLY at the very end.
40 from sklearn.model_selection import train_test_split
41 X_dev, X_holdout, y_dev, y_holdout = train_test_split(
42     X, y, test_size=0.15, random_state=42, stratify=y
43 )
44 print(f"      Dev set: {X_dev.shape[0]} samples | "
45       f"      Holdout: {X_holdout.shape[0]} samples")

```

```

46
47 #
48 # STEP 2: Hyperparameter tuning (inner CV loop)
49 #
50 print("\nStep 2  Hyperparameter Tuning (RandomizedSearchCV)...")
51
52 inner_cv = StratifiedKFold(n_splits=3, shuffle=True, random_state=
    42)
53
54 pipeline = Pipeline([
55     ('scaler', StandardScaler()),
56     ('clf', GradientBoostingClassifier(random_state=42))
57 ])
58
59 param_dist = {
60     'clf__n_estimators' : randint(50, 300),
61     'clf__max_depth'    : randint(2, 6),
62     'clf__learning_rate': loguniform(1e-2, 0.3),
63     'clf__subsample'    : uniform(0.5, 0.5),
64     'clf__min_samples_leaf': randint(5, 30),
65 }
66
67 search = RandomizedSearchCV(
68     pipeline, param_dist,
69     n_iter=30,
70     scoring='average_precision',
71     cv=inner_cv,
72     n_jobs=-1,
73     random_state=42,
74     refit=True
75 )
76 search.fit(X_dev, y_dev)
77 best_pipeline = search.best_estimator_
78
79 print(f"          Best tuning AP (inner CV): {search.best_score_:.4f}")
80 print(f"          Best params: n_est={search.best_params_['clf__n_
    estimators']}, "
81       f"depth={search.best_params_['clf__max_depth']}, "
82       f"lr={search.best_params_['clf__learning_rate']:.4f}")
83
84 #
85 # STEP 3: Unbiased performance estimation (outer CV loop)
86 #
87 print("\nStep 3  Outer CV: Unbiased Performance Estimation...")
88
89 outer_cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=0)
90
91 # For the outer loop we re-run the search inside each fold to
92 # avoid optimism, but for speed we demonstrate with the best
    pipeline.
93 # In practice, use a nested CV (GridSearchCV inside cross_validate).
94 cv_results = cross_validate(
95     best_pipeline, X_dev, y_dev,
96     cv=outer_cv,
97     scoring={
98         'ap'      : 'average_precision',
99         'auc'     : 'roc_auc',
100        'f1'      : 'f1',

```

```

101         'bal_acc': 'balanced_accuracy'
102     },
103     return_train_score=True,
104     n_jobs=-1
105 )
106
107 print(f"\n          {'Metric':<18} {'Mean':>8} {'Std':>8} {'Overfit Gap'
108       ':>12}")
109 print(f"          {'-'*48}")
110 for m in ['ap', 'auc', 'f1', 'bal_acc']:
111     train_mean = cv_results[f'train_{m}'].mean()
112     test_mean  = cv_results[f'test_{m}'].mean()
113     test_std   = cv_results[f'test_{m}'].std()
114     gap        = train_mean - test_mean
115     print(f"          {m:<18} {test_mean:>8.4f} {test_std:>8.4f} {gap:>
116           12.4f}")
117
118 #
119 # STEP 4: Threshold optimisation
120 #
121 print("\nStep 4  Threshold Optimisation...")
122
123 oof_proba = cross_val_predict(
124     best_pipeline, X_dev, y_dev, cv=outer_cv,
125     method='predict_proba'
126 )[:, 1]
127
128 prec, rec, thresholds = precision_recall_curve(y_dev, oof_proba)
129 # Use F2: recall matters 2x more than precision (missing positives
130 #         costly)
131 f2_scores = (5 * prec[:-1] * rec[:-1]) / (4*prec[:-1] + rec[:-1] + 1
132     e-9)
133 best_thresh_idx = np.argmax(f2_scores)
134 optimal_threshold = thresholds[best_thresh_idx]
135
136 print(f"          Optimal threshold (max F2): {optimal_threshold:.3f}")
137 print(f"          F2 at optimal: {f2_scores[best_thresh_idx]:.4f}")
138 print(f"          F2 at 0.50    : {f2_scores[np.argmin(np.abs(thresholds
139     -0.5))]:.4f}")
140
141 #
142 # STEP 5: Final holdout evaluation (one-time, no peeking before this
143 #         )
144 #
145 print("\nStep 5  Final Holdout Evaluation...")
146
147 best_pipeline.fit(X_dev, y_dev)
148 holdout_proba = best_pipeline.predict_proba(X_holdout)[:, 1]
149 holdout_labels = (holdout_proba >= optimal_threshold).astype(int)
150
151 print(f"\n          *** FINAL HOLDOUT RESULTS ***")
152 print(f"          AUC-ROC          : {roc_auc_score(y_holdout, holdout_
153     proba):.4f}")
154 print(f"          Average Precision: {average_precision_score(y_holdout
155     , holdout_proba):.4f}")
156 print(f"          Balanced Accuracy: {balanced_accuracy_score(y_holdout
157     , holdout_labels):.4f}")

```

```

149 print(f"          F1 Score          : {f1_score(y_holdout, holdout_
      labels):.4f}")
150 print()
151 print(classification_report(y_holdout, holdout_labels,
152                             target_names=['Negative', 'Positive']))
153
154 #
155 # STEP 6: SHAP interpretability on holdout set
156 #
157 print("Step 6  SHAP Global Interpretability...")
158
159 fitted_clf      = best_pipeline.named_steps['clf']
160 fitted_scaler   = best_pipeline.named_steps['scaler']
161 X_holdout_s     = fitted_scaler.transform(X_holdout)
162
163 explainer       = shap.TreeExplainer(fitted_clf)
164 shap_vals       = explainer.shap_values(X_holdout_s)
165
166 # Compute mean absolute SHAP and rank features
167 mean_abs_shap   = np.abs(shap_vals).mean(axis=0)
168 feat_importance = pd.Series(mean_abs_shap, index=feature_names)\
169                     .sort_values(ascending=False)
170
171 print("\n          Top 5 features by mean |SHAP|:")
172 for feat, val in feat_importance.head(5).items():
173     print(f"          {feat:<20}: {val:.5f}")
174
175 shap.summary_plot(shap_vals, X_holdout,
176                  feature_names=feature_names,
177                  plot_type='bar', show=False, max_display=10)
178 plt.tight_layout()
179 plt.savefig("final_shap_importance.png", dpi=150, bbox_inches='tight',)
180 plt.close()
181 print("\n          SHAP importance plot saved.")
182 print("\n" + "="*65)
183 print("  EVALUATION COMPLETE")
184 print("="*65)

```

6 The Evaluation Checklist

Before declaring a model production-ready, work through the following checklist. Each item addresses a common failure mode described in this tutorial.

Data Integrity

- ☐ All preprocessing fitted inside a Pipeline or on training data only.
- ☐ No future data in training splits (especially for time-series).
- ☐ Class distribution verified in each fold (use StratifiedKFold).
- ☐ Holdout test set segregated at the start and not inspected until the final evaluation.

Cross-Validation

- ☐ CV strategy matches deployment scenario (StratifiedKFold for classification, TimeSeriesSplit for temporal data).
- ☐ Both mean and standard deviation of CV scores reported.
- ☐ Training scores compared to validation scores to diagnose overfitting.
- ☐ Performance trend across time-series folds inspected for concept drift.

Hyperparameter Tuning

- ☐ Tuning performed on development set only; holdout never used during tuning.
- ☐ Nested cross-validation used if reporting unbiased performance alongside tuning.
- ☐ Convergence plot inspected to confirm budget was sufficient.

Metrics

- ☐ Metric chosen based on business cost structure, not convenience.
- ☐ For imbalanced datasets: accuracy not used alone; AP, F1, and balanced accuracy reported.
- ☐ Decision threshold optimised for the chosen business objective.
- ☐ Precision-recall curve inspected, not just the single-threshold result.

Interpretability

- ☐ SHAP global summary consistent with domain knowledge.
- ☐ At least one local explanation generated and reviewed for a representative positive prediction and a representative negative prediction.
- ☐ PDP or ICE curves used to verify expected monotonicity of key features.

- ☐ Any surprising model behaviour investigated for data errors or leakage before proceeding.

7 Regression Metrics: Measuring Continuous Predictions

7.1 The Problem with a Single Number

Classification has the luxury of a confusion matrix — a two-dimensional structure that separates types of errors. Regression collapses the comparison between predicted and actual values into a single error distribution. The danger is that different metrics collapse that distribution differently, each hiding a different kind of problem. A model can have a low RMSE but a terrible MAPE on small values, or a good R^2 but systematically biased predictions in one region. Understanding each metric's blind spots is what separates a careful practitioner from someone who reports the first number that looks good.

7.2 Core Regression Metrics

Let y_i be the true value and \hat{y}_i be the prediction, with n samples and mean true value \bar{y} .

Mean Absolute Error (MAE) is the average of the absolute differences. It is in the same unit as the target and gives every error equal weight:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Squared Error (MSE) squares each error before averaging, which means large errors are penalised disproportionately. This is appropriate when large errors are especially costly:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE) is the square root of MSE, restoring the original unit and making it directly comparable to MAE:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Mean Absolute Percentage Error (MAPE) expresses error as a percentage of the true value, making it scale-independent and useful for comparing models across datasets with different magnitudes. It is undefined or explosive when $y_i \approx 0$:

$$\text{MAPE} = \frac{100}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$

Symmetric MAPE (sMAPE) fixes the division-by-zero problem by dividing by the mean of the actual and predicted values:

$$\text{sMAPE} = \frac{100}{n} \sum_{i=1}^n \frac{2|y_i - \hat{y}_i|}{|y_i| + |\hat{y}_i|}$$

R^2 (Coefficient of Determination) measures the proportion of variance in the target explained by the model. A baseline model that always predicts \bar{y} achieves $R^2 = 0$; a perfect model achieves $R^2 = 1$. Negative values mean the model is worse than the mean:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Adjusted R^2 penalises for the number of features p , preventing R^2 from artificially increasing when irrelevant features are added:

$$\bar{R}^2 = 1 - (1 - R^2) \cdot \frac{n - 1}{n - p - 1}$$

Metric	Use when	Caution
MAE	Robust to outliers; median-optimising	Not differentiable at 0; less sensitive to large errors
RMSE	Large errors are disproportionately costly	Heavily influenced by outliers; unit-sensitive
MAPE	Comparing across scales; percentage makes sense	Blows up near zero; biased toward under-prediction
sMAPE	MAPE context but targets near zero exist	Can still behave oddly for very small values
R^2	Explaining variance; comparing to baseline	Scale-free but can be misleading with non-linear models
Adj. R^2	Feature selection; multivariate regression	Still a summary statistic; inspect residuals

Table 5: Regression metric selection guide.

7.3 Residual Analysis: Beyond a Single Number

No single metric captures the full picture. Residual analysis — examining the distribution and structure of errors — is the gold standard for diagnosing regression models.

Comprehensive Regression Evaluation with Residual Diagnostics

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib
4 matplotlib.use('Agg')
5 import matplotlib.pyplot as plt
6 from sklearn.datasets import fetch_california_housing

```



```

7 from sklearn.ensemble import GradientBoostingRegressor
8 from sklearn.linear_model import LinearRegression
9 from sklearn.model_selection import cross_val_predict, KFold, cross_
  validate
10 from sklearn.pipeline import Pipeline
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.metrics import (
13     mean_absolute_error, mean_squared_error,
14     r2_score, mean_absolute_percentage_error
15 )
16
17 # -----
18 # 1. California housing dataset: predict median house value
19 # Target is in units of $100,000. ~20,000 samples, 8 features.
20 # -----
21 housing = fetch_california_housing(as_frame=True)
22 X, y = housing.data, housing.target
23 print(f"Dataset: {X.shape[0]} samples, {X.shape[1]} features")
24 print(f"Target range: [{y.min():.2f}, {y.max():.2f}] "
25       f"Mean: {y.mean():.2f} Std: {y.std():.2f}")
26
27 # -----
28 # 2. Two models: linear baseline and gradient boosting
29 # -----
30 models = {
31     'Linear Regression': Pipeline([
32         ('scaler', StandardScaler()),
33         ('reg', LinearRegression())
34     ]),
35     'Gradient Boosting': Pipeline([
36         ('scaler', StandardScaler()),
37         ('reg', GradientBoostingRegressor(
38             n_estimators=200, learning_rate=0.05,
39             max_depth=5, random_state=42
40         ))
41     ])
42 }
43
44 cv = KFold(n_splits=5, shuffle=True, random_state=42)
45
46 def smape(y_true, y_pred):
47     """Symmetric Mean Absolute Percentage Error."""
48     denom = (np.abs(y_true) + np.abs(y_pred)) / 2
49     return np.mean(np.abs(y_true - y_pred) / denom) * 100
50
51 results = {}
52 for name, pipe in models.items():
53     # Out-of-fold predictions for unbiased metric computation
54     oof_pred = cross_val_predict(pipe, X, y, cv=cv, n_jobs=-1)
55
56     mae = mean_absolute_error(y, oof_pred)
57     rmse = np.sqrt(mean_squared_error(y, oof_pred))
58     mape = mean_absolute_percentage_error(y, oof_pred) * 100
59     smap = smape(y.values, oof_pred)
60     r2 = r2_score(y, oof_pred)
61
62     # Adjusted R^2
63     n, p = X.shape

```

```

64     adj_r2 = 1 - (1 - r2) * (n - 1) / (n - p - 1)
65
66     results[name] = {
67         'oof_pred': oof_pred,
68         'MAE': mae, 'RMSE': rmse,
69         'MAPE': mape, 'sMAPE': smape,
70         'R2': r2, 'Adj_R2': adj_r2
71     }
72
73     # -----
74     # 3. Metric comparison table
75     # -----
76     print(f"\n{'Metric':<12}", end="")
77     for name in models:
78         print(f"    {name:<22}", end="")
79     print()
80     print("-" * 60)
81     for metric in ['MAE', 'RMSE', 'MAPE', 'sMAPE', 'R2', 'Adj_R2']:
82         print(f"{metric:<12}", end="")
83         for name in models:
84             val = results[name][metric]
85             unit = "%" if metric in ['MAPE', 'sMAPE'] else ""
86             print(f"    {val:>8.4f}{unit:<14}", end="")
87         print()
88
89     # -----
90     # 4. Four-panel residual diagnostic plot (Gold Standard)
91     #     Panel A: Residuals vs Fitted    checks heteroscedasticity
92     #     Panel B: QQ plot                checks normality of errors
93     #     Panel C: Scale-Location        detects variance instability
94     #     Panel D: Error distribution     checks symmetry and tails
95     # -----
96     from scipy import stats
97
98     fig, axes = plt.subplots(2, 2, figsize=(13, 10))
99     best_name = 'Gradient Boosting'
100    oof = results[best_name]['oof_pred']
101    resid = y.values - oof
102    std_resid = resid / resid.std()
103
104    # Panel A: Residuals vs Fitted
105    axes[0,0].scatter(oof, resid, alpha=0.3, s=8, color='#0d47a1')
106    axes[0,0].axhline(0, color='red', linewidth=1.5, linestyle='--')
107    # Add LOWESS-style trend
108    from numpy.polynomial import polynomial as P
109    z = np.polyfit(oof, resid, 3)
110    p = np.poly1d(z)
111    x_line = np.linspace(oof.min(), oof.max(), 200)
112    axes[0,0].plot(x_line, p(x_line), color='orange', linewidth=2)
113    axes[0,0].set_xlabel("Fitted Values", fontsize=11)
114    axes[0,0].set_ylabel("Residuals", fontsize=11)
115    axes[0,0].set_title("Residuals vs Fitted\n(should be random around 0)",
116                        fontsize=11, fontweight='bold')
117    axes[0,0].grid(alpha=0.3)
118
119    # Panel B: Normal Q-Q Plot

```

```

120 (osm, osr), (slope, intercept, r) = stats.probplot(std_resid, dist='
    norm')
121 axes[0,1].scatter(osm, osr, alpha=0.4, s=8, color='#0d47a1')
122 axes[0,1].plot(osm, slope*np.array(osm)+intercept, 'r--', linewidth=
    1.5)
123 axes[0,1].set_xlabel("Theoretical Quantiles", fontsize=11)
124 axes[0,1].set_ylabel("Sample Quantiles", fontsize=11)
125 axes[0,1].set_title("Normal Q-Q Plot\n(deviations = non-normal
    errors)",
    fontsize=11, fontweight='bold')
126
127 axes[0,1].grid(alpha=0.3)
128
129 # Panel C: Scale-Location (sqrt of abs residuals vs fitted)
130 axes[1,0].scatter(oof, np.sqrt(np.abs(std_resid)), alpha=0.3, s=8,
    color='#00838f')
131
132 axes[1,0].set_xlabel("Fitted Values", fontsize=11)
133 axes[1,0].set_ylabel(r"$\sqrt{|\text{Std. Residuals}|}$", fontsize=
    11)
134 axes[1,0].set_title("Scale-Location\n(flat line = homoscedastic)",
    fontsize=11, fontweight='bold')
135
136 axes[1,0].grid(alpha=0.3)
137
138 # Panel D: Error distribution
139 axes[1,1].hist(resid, bins=60, color='#0d47a1', alpha=0.75,
    edgecolor='white',
140     linewidth=0.3, density=True)
141 xr = np.linspace(resid.min(), resid.max(), 300)
142 axes[1,1].plot(xr, stats.norm.pdf(xr, resid.mean(), resid.std()),
    'r--', linewidth=2, label='Normal fit')
143
144 axes[1,1].set_xlabel("Residual", fontsize=11)
145 axes[1,1].set_ylabel("Density", fontsize=11)
146 axes[1,1].set_title(f"Error Distribution\n"
    f"Skew={stats.skew(resid):.3f} "
    f"Kurt={stats.kurtosis(resid):.3f}",
    fontsize=11, fontweight='bold')
147
148
149 axes[1,1].legend(fontsize=9)
150 axes[1,1].grid(alpha=0.3)
151
152
153 plt.suptitle(f"Residual Diagnostics {best_name}",
    fontsize=14, fontweight='bold', y=1.01)
154
155 plt.tight_layout()
156 plt.savefig("residual_diagnostics.png", dpi=150, bbox_inches='tight'
    )
157 plt.close()
158 print("\nResidual diagnostic plots saved.")
159
160 # -----
161 # 5. Prediction interval width analysis
162 # Pinball / quantile loss: for each quantile tau, a model that
163 # outputs the tau-quantile should have tau fraction of actuals
    below.
164 # -----
165 from sklearn.ensemble import GradientBoostingRegressor as GBR
166
167 quantiles = [0.1, 0.5, 0.9]
168 coverage_results = {}
169
170 for q in quantiles:

```

```

171     q_model = Pipeline([
172         ('scaler', StandardScaler()),
173         ('reg', GBR(loss='quantile', alpha=q,
174             n_estimators=150, max_depth=4, random_state=42))
175     ])
176     q_pred = cross_val_predict(q_model, X, y, cv=cv, n_jobs=-1)
177     coverage = np.mean(y.values <= q_pred)
178     coverage_results[q] = coverage
179
180     print("\nQuantile calibration (target coverage vs actual coverage):"
181         )
182     for q, cov in coverage_results.items():
183         print(f"    Q{int(q*100):>2}: target={q:.2f}    actual={cov:.3f}    "
184             f"{ 'OK' if abs(q - cov) < 0.03 else 'MISCALIBRATED' }")

```

The **four-panel residual plot** is the equivalent of the confusion matrix for regression. The *Residuals vs Fitted* panel should show a random cloud centred at zero — any curve or funnel shape indicates the model is missing systematic structure. The *Q-Q plot* should follow the diagonal line — heavy tails mean the model is underpredicting the frequency of large errors. The *Scale-Location* plot should be flat — a rising trend means variance grows with the prediction (heteroscedasticity), which invalidates confidence intervals.

8 Model Calibration: When Probabilities Must Mean Something

8.1 The Calibration Problem

A classifier's predicted probability is not just a ranking device — in many business contexts it is used directly as a probability estimate. A credit scoring model that outputs 0.30 for a customer is implicitly claiming that 30% of customers with that profile will default. If in reality 60% of those customers default, the model is severely *miscalibrated*, and any downstream decision made using the raw probability (e.g., expected loss calculation, premium pricing) will be systematically wrong.

Calibration is the statistical property that a model's predicted probability of p should match the observed frequency of the event across all predictions in the neighbourhood of p .

Formally, a model is perfectly calibrated if:

$$P(Y = 1 \mid \hat{p} = p) = p \quad \forall p \in [0, 1]$$

8.2 The Reliability Diagram

The *reliability diagram* (also called the calibration curve) is the primary visual tool. Predictions are grouped into bins by their predicted probability; for each bin, the fraction of actual positive outcomes is plotted against the mean predicted probability. A perfectly calibrated model lies on the $y = x$ diagonal.

8.3 The Brier Score

The *Brier score* is the mean squared error of probability predictions:

$$\text{BS} = \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i)^2$$

It ranges from 0 (perfect) to 1 (perfectly wrong). It penalises confident wrong predictions heavily and decomposes into:

$$\text{BS} = \underbrace{\text{Reliability}}_{\text{calibration error}} - \underbrace{\text{Resolution}}_{\text{sharpness}} + \underbrace{\text{Uncertainty}}_{\text{base rate variance}}$$

A model can have a low Brier score either by being well calibrated or by being sharp (confident). Reporting both Brier score and the reliability diagram separates the two contributions.

Calibration Analysis, Platt Scaling, and Isotonic Regression

```
1 import numpy as np
2 import matplotlib
3 matplotlib.use('Agg')
```

```

4 import matplotlib.pyplot as plt
5 from sklearn.datasets import make_classification
6 from sklearn.ensemble import GradientBoostingClassifier,
  RandomForestClassifier
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.calibration import CalibratedClassifierCV, calibration_
  curve
9 from sklearn.model_selection import train_test_split, cross_val_
  predict, StratifiedKFold
10 from sklearn.pipeline import Pipeline
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.metrics import brier_score_loss
13
14 # -----
15 # 1. Dataset and train/test split
16 # -----
17 X, y = make_classification(
18     n_samples=5000, n_features=20, n_informative=10,
19     weights=[0.75, 0.25], random_state=42
20 )
21 X_train, X_test, y_train, y_test = train_test_split(
22     X, y, test_size=0.3, stratify=y, random_state=42
23 )
24
25 # -----
26 # 2. Gradient Boosting is often poorly calibrated because it was
27 #    optimised for ranking (log-loss), not for probability accuracy.
28 #    Random Forest is typically overconfident (predicted probs
29 #    cluster toward 0 and 1, not matching actual frequencies).
30 # -----
31 models = {
32     'Gradient Boosting (raw)': Pipeline([
33         ('scaler', StandardScaler()),
34         ('clf', GradientBoostingClassifier(
35             n_estimators=200, max_depth=4, random_state=42
36         ))
37     ]),
38     'Random Forest (raw)': Pipeline([
39         ('scaler', StandardScaler()),
40         ('clf', RandomForestClassifier(
41             n_estimators=200, random_state=42
42         ))
43     ]),
44 }
45
46 # -----
47 # 3. Calibrated versions
48 #    Platt Scaling (method='sigmoid'): fits a logistic regression
49 #    on top of the model's raw scores. Works well when the
50 #    calibration curve is S-shaped.
51 #
52 #    Isotonic Regression (method='isotonic'): fits a non-parametric
53 #    monotonic function. More flexible but needs more data to avoid
54 #    overfitting. Use when the calibration error is non-monotonic.
55 #
56 #    cv='prefit': the base model is already fitted; only the
57 #    calibration layer is fitted on the calibration set.
58 # -----

```

```

59 calibrated_models = {}
60 for name, pipe in models.items():
61     pipe.fit(X_train, y_train)
62     # Platt scaling
63     calib_platt = CalibratedClassifierCV(pipe, method='sigmoid', cv=
        'prefit')
64     calib_platt.fit(X_test[:len(X_test)//2], y_test[:len(y_test)//2]
        )
65     calibrated_models[name + ' + Platt'] = calib_platt
66
67     # Isotonic regression
68     calib_iso = CalibratedClassifierCV(pipe, method='isotonic', cv='
        prefit')
69     calib_iso.fit(X_test[:len(X_test)//2], y_test[:len(y_test)//2])
70     calibrated_models[name + ' + Isotonic'] = calib_iso
71
72 # Use the remaining half of test set for evaluation
73 X_eval, y_eval = X_test[len(X_test)//2:], y_test[len(y_test)//2:]
74
75 # -----
76 # 4. Compute calibration curves and Brier scores
77 # -----
78 all_models = {**models, **calibrated_models}
79 eval_results = {}
80
81 for name, model in all_models.items():
82     try:
83         proba = model.predict_proba(X_eval)[: , 1]
84     except Exception:
85         proba = model.predict_proba(X_eval)[: , 1]
86
87     frac_pos, mean_pred = calibration_curve(
88         y_eval, proba, n_bins=10, strategy='uniform'
89     )
90     brier = brier_score_loss(y_eval, proba)
91     eval_results[name] = {
92         'proba': proba, 'frac_pos': frac_pos,
93         'mean_pred': mean_pred, 'brier': brier
94     }
95
96 # -----
97 # 5. Reliability diagram: 2 x 2 grid comparing raw vs calibrated
98 # -----
99 fig, axes = plt.subplots(1, 2, figsize=(13, 6))
100
101 colors_raw = ['#0d47a1', '#00838f']
102 colors_cal = ['#1565c0', '#e65100', '#1b5e20', '#880e4f']
103
104 for ax, base_name, color_raw in zip(
105     axes, ['Gradient Boosting', 'Random Forest'], colors_raw):
106
107     ax.plot([0, 1], [0, 1], 'k--', linewidth=1.2, label='Perfect
        calibration')
108
109     # Raw model
110     r = eval_results[base_name + ' (raw)']
111     ax.plot(r['mean_pred'], r['frac_pos'], 's-', color=color_raw,
112         linewidth=2, markersize=6,

```

```

113         label=f"Raw (Brier={r['brier']:.4f})")
114
115     # Platt
116     rp = eval_results[base_name + ' (raw) + Platt']
117     ax.plot(rp['mean_pred'], rp['frac_pos'], 'o-', color='#e65100',
118            linewidth=2, markersize=6,
119            label=f"Platt (Brier={rp['brier']:.4f})")
120
121     # Isotonic
122     ri = eval_results[base_name + ' (raw) + Isotonic']
123     ax.plot(ri['mean_pred'], ri['frac_pos'], '^-', color='#1b5e20',
124            linewidth=2, markersize=6,
125            label=f"Isotonic (Brier={ri['brier']:.4f})")
126
127     ax.set_xlabel("Mean Predicted Probability", fontsize=11)
128     ax.set_ylabel("Fraction of Positives", fontsize=11)
129     ax.set_title(f"Reliability Diagram {base_name.split(' ')[0]}",
130                fontsize=12, fontweight='bold')
131     ax.legend(fontsize=8.5)
132     ax.grid(alpha=0.3)
133     ax.set_xlim(0, 1); ax.set_ylim(0, 1)
134
135 plt.tight_layout()
136 plt.savefig("calibration_curves.png", dpi=150, bbox_inches='tight')
137 plt.close()
138
139 # -----
140 # 6. Summary: Brier scores
141 # -----
142 print("Brier Score Summary (lower = better):")
143 for name, r in eval_results.items():
144     print(f" {name:<42}: {r['brier']:.5f}")
145 print("\nNote: calibration only improves probability accuracy, not")
146 print("discrimination (AUC). Always report both AUC and Brier score.")

```

When calibration matters most: any time probabilities feed into a downstream system that treats them as actual probabilities — expected value calculations, risk-adjusted pricing, multi-model ensembles that average probabilities, or Bayesian updating. When the model output is used only as a ranking (e.g., top-K recommendations, ranked lists), calibration is irrelevant and AUC is the right metric.

Calibration and discrimination (AUC) are orthogonal properties. You can have a perfectly discriminating model with terrible calibration (a monotonic transformation of the probabilities preserves AUC but destroys calibration). Always report both. Use AUC to assess whether the model separates classes; use the Brier score and reliability diagram to assess whether the probabilities are trustworthy.

9 Statistical Significance Testing for Model Comparison

9.1 The Core Question: Is Model A Actually Better?

When Model A achieves an AUC of 0.872 and Model B achieves 0.861 on the same cross-validation folds, the natural question is: is this difference real, or is it sampling noise? Without a statistical test, you cannot answer that question. Choosing the better-looking model on a single evaluation is like flipping a coin twice and concluding it is biased toward heads.

The challenge in model comparison is that the observations (CV fold scores) are *not independent* — they share training data. Standard two-sample t-tests assume independence and are anti-conservative (too likely to declare a significant difference) in this setting.

9.2 McNemar's Test: Comparing Classifiers on the Same Test Set

McNemar's test compares two classifiers on the same set of test examples. It asks: of the examples where the two models disagree, does one model get more right than the other? It operates on the contingency table of per-sample correct/incorrect decisions and requires no distributional assumptions.

$$\chi^2 = \frac{(n_{01} - n_{10})^2}{n_{01} + n_{10}}$$

where n_{01} is the number of examples Model A gets wrong but Model B gets right, and n_{10} is the reverse.

9.3 The Corrected Resampled t-Test

For comparing models across cross-validation folds, Nadeau and Bengio (2003) proposed a correction to the paired t-test that accounts for the overlap in training data between folds. For k -fold CV with n total samples and $n_{\text{test}} = n/k$ test samples per fold:

$$t = \frac{\bar{d}}{\sqrt{\left(\frac{1}{k} + \frac{n_{\text{test}}}{n_{\text{train}}}\right) s_d^2}}$$

where \bar{d} is the mean per-fold score difference and s_d^2 is the variance of those differences.

Statistical Model Comparison: McNemar and Corrected t-Test

```
1 import numpy as np
2 import pandas as pd
3 from scipy import stats
4 from scipy.stats import chi2
5 from sklearn.datasets import make_classification
6 from sklearn.ensemble import GradientBoostingClassifier,
```

```

    RandomForestClassifier
7  from sklearn.linear_model import LogisticRegression
8  from sklearn.model_selection import StratifiedKFold, cross_val_
   predict, cross_val_score
9  from sklearn.pipeline import Pipeline
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.metrics import roc_auc_score
12 from itertools import combinations
13
14 X, y = make_classification(
15     n_samples=3000, n_features=20, n_informative=10,
16     weights=[0.7, 0.3], random_state=42
17 )
18
19 cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
20
21 models = {
22     'Logistic Regression': Pipeline([
23         ('sc', StandardScaler()),
24         ('clf', LogisticRegression(C=0.1, max_iter=1000))
25     ]),
26     'Random Forest': Pipeline([
27         ('sc', StandardScaler()),
28         ('clf', RandomForestClassifier(n_estimators=200, random_
   state=42))
29     ]),
30     'Gradient Boosting': Pipeline([
31         ('sc', StandardScaler()),
32         ('clf', GradientBoostingClassifier(
33             n_estimators=200, learning_rate=0.05,
34             max_depth=4, random_state=42
35         ))
36     ]),
37 }
38
39 # -----
40 # 1. Collect per-fold AUC scores and per-sample hard predictions
41 # -----
42 fold_scores = {}
43 oof_labels = {}
44 n = len(y)
45 n_splits = cv.get_n_splits()
46 n_test_per_fold = n // n_splits
47
48 for name, pipe in models.items():
49     fold_aucs = cross_val_score(pipe, X, y, cv=cv, scoring='roc_auc',
   , n_jobs=-1)
50     oof_preds = cross_val_predict(pipe, X, y, cv=cv, n_jobs=-1) #
   hard labels
51     fold_scores[name] = fold_aucs
52     oof_labels[name] = oof_preds
53     print(f"{name:<25}: AUC={fold_aucs.mean():.4f} +/- {fold_aucs.
   std():.4f}")
54
55 # -----
56 # 2. Corrected Resampled t-Test (Nadeau-Bengio)
57 # Compares two models using their per-fold score differences.
58 # The correction term accounts for the non-independence of CV

```

```

folds.
59 # -----
60 def corrected_resampled_t_test(scores_a, scores_b, n_total, n_test):
61     """
62     Nadeau-Bengio corrected paired t-test for cross-validation
        comparison.
63
64     Parameters
65     -----
66     scores_a, scores_b : arrays of per-fold scores
67     n_total : total number of samples in dataset
68     n_test : number of test samples per fold
69
70     Returns
71     -----
72     t_stat, p_value
73     """
74     k = len(scores_a)
75     d = scores_a - scores_b # per-fold differences
76     d_bar = d.mean()
77     s2_d = d.var(ddof=1) # sample variance
78
79     n_train = n_total - n_test
80     # Correction factor: 1/k + n_test/n_train
81     correction = (1.0 / k) + (n_test / n_train)
82     se = np.sqrt(correction * s2_d)
83
84     if se == 0:
85         return 0.0, 1.0
86     t_stat = d_bar / se
87     p_val = 2 * stats.t.sf(np.abs(t_stat), df=k - 1)
88     return t_stat, p_val
89
90 print("\n--- Corrected Resampled t-Test (Nadeau-Bengio) ---")
91 print(f"{'Comparison':<45} {'t-stat':>8} {'p-value':>10} {'Significant':>12}")
92 print("-" * 78)
93
94 model_names = list(models.keys())
95 for a, b in combinations(model_names, 2):
96     t, p = corrected_resampled_t_test(
97         fold_scores[a], fold_scores[b],
98         n_total=n, n_test=n_test_per_fold
99     )
100     sig = "YES (p<0.05)" if p < 0.05 else "no"
101     print(f"{a} vs {b:<20} {t:>8.3f} {p:>10.4f} {sig:>12}")
102
103 # -----
104 # 3. McNemar's Test: comparing hard-label decisions on same samples
105 # Only meaningful for samples where the two models disagree.
106 # -----
107 def mcnemar_test(labels_a, labels_b, y_true):
108     """
109     Test whether two classifiers differ significantly in their
        errors.
110     n01: A wrong, B right
111     n10: A right, B wrong
112     Returns chi-squared statistic and p-value (with continuity

```

```

        correction).
113     """
114     n01 = np.sum((labels_a != y_true) & (labels_b == y_true))
115     n10 = np.sum((labels_a == y_true) & (labels_b != y_true))
116     # Continuity-corrected McNemar
117     if n01 + n10 == 0:
118         return 0.0, 1.0
119     chi2_stat = (abs(n01 - n10) - 1)**2 / (n01 + n10)
120     p_val = 1 - chi2.cdf(chi2_stat, df=1)
121     return chi2_stat, p_val
122
123 print("\n--- McNemar's Test (on out-of-fold hard labels) ---")
124 print(f"{'Comparison':<45} {'chi2':>8} {'p-value':>10} "
125       f"{'n01':>6} {'n10':>6} {'Significant':>12}")
126 print("-" * 95)
127
128 for a, b in combinations(model_names, 2):
129     chi2_val, p = mcnemar_test(oof_labels[a], oof_labels[b], y)
130     n01 = np.sum((oof_labels[a] != y) & (oof_labels[b] == y))
131     n10 = np.sum((oof_labels[a] == y) & (oof_labels[b] != y))
132     sig = "YES (p<0.05)" if p < 0.05 else "no"
133     print(f"{a} vs {b:<20} {chi2_val:>8.3f} {p:>10.4f} "
134           f"{n01:>6} {n10:>6} {sig:>12}")
135
136 # -----
137 # 4. Multiple comparisons correction (Bonferroni)
138 # When testing k pairs, the chance of a false positive rises.
139 # Bonferroni adjusts the significance threshold to alpha/k.
140 # -----
141 n_comparisons = len(list(combinations(model_names, 2)))
142 alpha_bonf = 0.05 / n_comparisons
143 print(f"\nBonferroni-corrected significance threshold for {n_
144       comparisons} "
145       f"comparisons: {alpha_bonf:.4f}")

```

If you test 20 model variants and declare significance at $p < 0.05$, you expect roughly one false positive by chance alone. Always apply a correction: Bonferroni (divide threshold by number of tests, conservative) or Benjamini-Hochberg (controls the false discovery rate, less conservative). Report corrected p-values whenever comparing more than two models.

10 Bias & Fairness Evaluation

10.1 Why Fairness Is a Measurement Problem

Fairness is not a single property — it is a family of mathematical definitions that are often mutually contradictory. The practitioner’s job is to: (1) understand which definition of fairness is appropriate for the business and regulatory context, (2) measure it, and (3) make explicit trade-offs when definitions conflict.

Before any fairness analysis, a critical distinction must be made. A *protected attribute* is a feature (race, gender, age, disability) that must not be the basis for a decision. Such features may be excluded from the model, but they must still be used to *evaluate* whether the model produces disparate outcomes.

10.2 Core Fairness Definitions

Demographic Parity (also: statistical parity): the model’s positive prediction rate should be equal across groups. Used in hiring and lending where equal rates of opportunity are required.

$$P(\hat{Y} = 1 \mid A = 0) = P(\hat{Y} = 1 \mid A = 1)$$

Equalized Odds: both the true positive rate (recall) and false positive rate must be equal across groups. Used in criminal justice and medical screening where errors of both types must not be concentrated in one group.

$$P(\hat{Y} = 1 \mid Y = y, A = 0) = P(\hat{Y} = 1 \mid Y = y, A = 1) \quad \forall y \in \{0, 1\}$$

Equal Opportunity: a relaxation of equalized odds that only requires equal recall (true positive rate). Appropriate when false negatives are the primary concern.

Disparate Impact Ratio: a ratio metric used in US employment law (the “80% rule”). If the positive rate for the disadvantaged group is less than 80% of the advantaged group’s rate, disparate impact is indicated.

$$\text{DIR} = \frac{P(\hat{Y} = 1 \mid A = 1)}{P(\hat{Y} = 1 \mid A = 0)} \quad (\text{should be } \geq 0.8)$$

Fairness Audit: Measuring and Visualising Bias

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib
4 matplotlib.use('Agg')
5 import matplotlib.pyplot as plt
6 from sklearn.datasets import make_classification
7 from sklearn.ensemble import GradientBoostingClassifier
8 from sklearn.model_selection import train_test_split,
   StratifiedKFold, cross_val_predict
9 from sklearn.pipeline import Pipeline
```

```

10 from sklearn.preprocessing import StandardScaler
11 from sklearn.metrics import (confusion_matrix, roc_auc_score,
12                               precision_score, recall_score)
13
14 np.random.seed(42)
15 n = 4000
16
17 # -----
18 # 1. Simulated loan application dataset
19 #   Group A=0: majority group (60%)
20 #   Group A=1: minority group (40%)
21 #   Introduce historical bias: minority group has artificially
22 #   lower income in training data (reflecting past discrimination).
23 # -----
24 group = np.random.binomial(1, 0.4, n) # 0=majority, 1=minority
25
26 income = np.where(group == 0,
27                   np.random.normal(65, 15, n), # majority: mean
28                   np.random.normal(52, 14, n)) # minority: mean
29
30 credit = np.random.normal(680, 80, n) - 30 * group
31 dti     = np.random.beta(2, 5, n) + 0.05 * group
32 history = np.random.gamma(4, 2, n) - 1.5 * group
33
34 X = np.column_stack([income, credit, dti, history,
35                     np.random.randn(n, 6)]) # 6 noise features
36 feature_names = ['income', 'credit_score', 'dti', 'credit_history',
37                 *[f'noise_{i}' for i in range(6)]]
38
39 # True default driven by financial fundamentals (fair ground truth)
40 logit = -0.04*income + 0.005*(750-credit) + 2*dti - 0.05*history
41 y      = (1/(1+np.exp(-logit + np.random.randn(n)*0.5))) > 0.5).astype
42         (int)
43
44 print(f"Default rate majority: {y[group==0].mean():.3f} "
45       f"minority: {y[group==1].mean():.3f}")
46
47 # -----
48 # 2. Train a model WITHOUT group as a feature (group-blind)
49 #   Fairness analysis uses group ONLY in evaluation, not training.
50 # -----
51 X_train, X_test, y_train, y_test, g_train, g_test = train_test_split
52     (
53     X, y, group, test_size=0.3, stratify=y, random_state=42
54 )
55
56 pipeline = Pipeline([
57     ('scaler', StandardScaler()),
58     ('clf', GradientBoostingClassifier(
59         n_estimators=200, max_depth=4,
60         learning_rate=0.05, random_state=42
61     ))
62 ])
63 pipeline.fit(X_train, y_train)
64 proba_test = pipeline.predict_proba(X_test)[:, 1]
65 pred_test  = pipeline.predict(X_test)

```

```

64
65 # -----
66 # 3. Per-group metrics
67 # -----
68 def group_metrics(y_true, y_pred, y_proba, groups):
69     """Compute fairness-relevant metrics per group."""
70     rows = []
71     for g_val in sorted(np.unique(groups)):
72         mask = groups == g_val
73         yt = y_true[mask]
74         yp = y_pred[mask]
75         ypr = y_proba[mask]
76         cm = confusion_matrix(yt, yp, labels=[0,1])
77         tn, fp, fn, tp = cm.ravel()
78         rows.append({
79             'group'           : g_val,
80             'n'               : mask.sum(),
81             'pos_rate'        : yp.mean(),          # demographic parity
82             'metric'          :
83             'tpr'             : tp/(tp+fn+1e-9),    # recall / equal
84             'opportunity'      :
85             'fpr'             : fp/(fp+tn+1e-9),    # false positive
86             'rate'            :
87             'precision'       : precision_score(yt, yp, zero_division=
88             0),
89             'auc'             : roc_auc_score(yt, ypr),
90             'false_neg_rate' : fn/(fn+tp+1e-9),
91         })
92     return pd.DataFrame(rows)
93
94 gm = group_metrics(y_test, pred_test, proba_test, g_test)
95 print("\n--- Per-Group Metrics ---")
96 print(gm.set_index('group').round(4).to_string())
97
98 # -----
99 # 4. Fairness metrics derived from group metrics
100 # -----
101 g0 = gm[gm['group']==0].iloc[0]
102 g1 = gm[gm['group']==1].iloc[0]
103
104 print("\n--- Fairness Metrics ---")
105 dir_ratio = g1['pos_rate'] / (g0['pos_rate'] + 1e-9)
106 demo_gap = abs(g0['pos_rate'] - g1['pos_rate'])
107 tpr_gap = abs(g0['tpr'] - g1['tpr'])
108 fpr_gap = abs(g0['fpr'] - g1['fpr'])
109 eq_odds = max(tpr_gap, fpr_gap)
110
111 print(f" Disparate Impact Ratio (>=0.8 required) : {dir_ratio:.4f}
112 "
113       f"{'PASS' if dir_ratio >= 0.8 else 'FAIL'}")
114 print(f" Demographic Parity Gap (|g0 - g1|) : {demo_gap:.4f} "
115       f"{'OK' if demo_gap < 0.05 else 'CONCERN'}")
116 print(f" Equal Opportunity Gap (TPR diff) : {tpr_gap:.4f} "
117       f"{'OK' if tpr_gap < 0.05 else 'CONCERN'}")
118 print(f" Equalized Odds Violation (max gap) : {eq_odds:.4f} "
119       f"{'OK' if eq_odds < 0.05 else 'CONCERN'}")
120
121 # -----

```

```

117 # 5. Threshold adjustment for equalized odds
118 # Use a different decision threshold per group to equalise TPR.
119 # This is the most common post-processing fairness intervention.
120 # -----
121 from sklearn.metrics import roc_curve
122
123 thresholds_explored = np.linspace(0.1, 0.9, 200)
124
125 # Find threshold per group that achieves TPR close to target
126 target_tpr = 0.70 # desired recall for positive class in both
127 # groups
128
129 best_thresh = {}
130 for g_val in [0, 1]:
131     mask = g_test == g_val
132     yt = y_test[mask]
133     ypr = proba_test[mask]
134     tprs = [recall_score(yt, (ypr >= t).astype(int), zero_division=
135         0)
136             for t in thresholds_explored]
137     best_idx = np.argmin(np.abs(np.array(tprs) - target_tpr))
138     best_thresh[g_val] = thresholds_explored[best_idx]
139
140 print(f"\n Threshold adjustment for TPR={target_tpr}:")
141 print(f" Group 0 threshold: {best_thresh[0]:.3f}")
142 print(f" Group 1 threshold: {best_thresh[1]:.3f}")
143
144 adj_pred = np.where(
145     g_test == 0,
146     (proba_test >= best_thresh[0]).astype(int),
147     (proba_test >= best_thresh[1]).astype(int)
148 )
149
150 gm_adj = group_metrics(y_test, adj_pred, proba_test, g_test)
151 print("\n--- After Threshold Adjustment ---")
152 print(gm_adj.set_index('group')[['pos_rate', 'tpr', 'fpr', 'precision']
153     ].round(4))
154
155 # -----
156 # 6. Visualise per-group metric comparison
157 # -----
158
159 fig, axes = plt.subplots(1, 2, figsize=(12, 5))
160
161 metrics_to_plot = ['pos_rate', 'tpr', 'fpr', 'precision', 'auc']
162 labels = ['Pos. Rate', 'TPR', 'FPR', 'Precision', 'AUC']
163
164 for ax, (gm_df, title) in zip(axes, [
165     (gm, 'Before Adjustment'),
166     (gm_adj, 'After Threshold Adjustment')
167 ]):
168     x = np.arange(len(metrics_to_plot))
169     w = 0.35
170     vals0 = [gm_df[gm_df['group']==0].iloc[0][m] for m in metrics_to
171         _plot]
172     vals1 = [gm_df[gm_df['group']==1].iloc[0][m] for m in metrics_to
173         _plot]
174     ax.bar(x - w/2, vals0, w, label='Group 0 (majority)',

```



```

170         color='#0d47a1', alpha=0.85)
171     ax.bar(x + w/2, vals1, w, label='Group 1 (minority)',
172           color='#e65100', alpha=0.85)
173     ax.set_xticks(x); ax.set_xticklabels(labels, fontsize=10)
174     ax.set_ylim(0, 1.1)
175     ax.axhline(0.8, color='red', linestyle='--', linewidth=1,
176              label='80% rule threshold')
177     ax.set_title(title, fontsize=12, fontweight='bold')
178     ax.legend(fontsize=8); ax.grid(axis='y', alpha=0.3)
179
180 plt.suptitle("Fairness Metrics Before and After Threshold Adjustment",
181             ,
182             fontsize=13, fontweight='bold')
183 plt.tight_layout()
184 plt.savefig("fairness_metrics.png", dpi=150, bbox_inches='tight')
185 plt.close()
186 print("\nFairness visualisation saved.")

```

Chouldechova (2017) and Kleinberg et al. (2017) independently proved that demographic parity, equalized odds, and predictive parity (precision equality) *cannot all be satisfied simultaneously* unless base rates are equal across groups. Fairness analysis requires an explicit choice of which criterion to prioritise, determined by the ethical and legal context of the application, not by the algorithm.

11 Production Monitoring & Drift Detection

11.1 The Life Cycle Does Not End at Deployment

A model that performs well on the test set at deployment time may degrade over weeks or months as the world changes. This degradation has two root causes.

Data drift (or covariate shift) occurs when the distribution of input features $P(X)$ changes while the relationship $P(Y | X)$ remains stable. Example: a pandemic shifts the income distribution of loan applicants, but the relationship between income and default probability is unchanged.

Concept drift occurs when $P(Y | X)$ itself changes — the fundamental relationship between features and target has shifted. Example: a new regulation changes what constitutes fraud, so past labelling conventions no longer apply.

Both require monitoring, but they call for different responses: data drift may be handled by retraining on new data; concept drift may require redesigning the feature set or the label definition.

11.2 Population Stability Index (PSI)

The **Population Stability Index** is the most widely used metric for monitoring input drift in production. It was developed by credit risk practitioners and measures the divergence between a reference distribution and a current distribution.

$$\text{PSI} = \sum_{b=1}^B (q_b - p_b) \ln \left(\frac{q_b}{p_b} \right)$$

where p_b is the proportion of reference data in bin b and q_b is the corresponding proportion in the current data. PSI is a symmetrised version of the KL divergence. The standard thresholds are: $\text{PSI} < 0.1$ (stable), $0.1 \leq \text{PSI} < 0.2$ (minor shift, monitor closely), $\text{PSI} \geq 0.2$ (significant shift, investigate and consider retraining).

PSI, KS Drift Detection, and Performance Monitoring

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib
4 matplotlib.use('Agg')
5 import matplotlib.pyplot as plt
6 from scipy import stats
7 from sklearn.datasets import make_classification
8 from sklearn.ensemble import GradientBoostingClassifier
9 from sklearn.model_selection import train_test_split
10 from sklearn.pipeline import Pipeline
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.metrics import roc_auc_score
13
14 # -----
15 # 1. Simulate production lifecycle
16 #    Month 0-2: stable data (training distribution)

```

```

17 #     Month 3-5: gradual covariate shift (data drift)
18 #     Month 6-8: concept drift (relationship changes)
19 # -----
20 np.random.seed(42)
21
22 def generate_batch(n, drift_factor=0.0, concept_drift=False):
23     """Generate a batch of data with optional drift."""
24     X = np.random.randn(n, 10)
25     X[:, 0] += drift_factor * 2      # feature 0 shifts with drift
26     X[:, 1] += drift_factor * 1.5    # feature 1 also shifts
27
28     if not concept_drift:
29         # Original relationship: feature 0 and 2 predict y
30         logit = 1.5*X[:,0] + X[:,2] + 0.3*np.random.randn(n)
31     else:
32         # Concept drift: feature 0 now has OPPOSITE sign
33         logit = -1.5*X[:,0] + X[:,2] + 0.3*np.random.randn(n)
34
35     y = (1/(1 + np.exp(-logit)) > 0.5).astype(int)
36     return X, y
37
38 # Training data (reference)
39 X_train, y_train = generate_batch(3000, drift_factor=0.0)
40
41 # Production batches (one per month)
42 months = list(range(9))
43 monthly_data = {
44     m: generate_batch(
45         300,
46         drift_factor=max(0, (m-2)*0.4),      # drift starts month
47         concept_drift=(m >= 6)                # concept drift
48         from month 6
49     )
50     for m in months
51 }
52 # -----
53 # 2. Train model on reference data
54 # -----
55 pipeline = Pipeline([
56     ('scaler', StandardScaler()),
57     ('clf', GradientBoostingClassifier(n_estimators=150, random_
58         state=42))
59 ])
60 pipeline.fit(X_train, y_train)
61 # -----
62 # 3. PSI function
63 # -----
64 def compute_psi(reference, current, n_bins=10):
65     """
66     Compute Population Stability Index between two samples.
67     Uses percentile-based bins from the reference distribution.
68     """
69     breakpoints = np.percentile(reference, np.linspace(0, 100, n_
70         bins+1))
71     breakpoints = np.unique(breakpoints) # remove duplicates

```

```

71
72     ref_counts = np.histogram(reference, bins=breakpoints)[0]
73     cur_counts = np.histogram(current, bins=breakpoints)[0]
74
75     # Convert to proportions, add small epsilon to avoid log(0)
76     eps = 1e-6
77     ref_pct = ref_counts / len(reference) + eps
78     cur_pct = cur_counts / len(current) + eps
79
80     psi = np.sum((cur_pct - ref_pct) * np.log(cur_pct / ref_pct))
81     return psi
82
83 # -----
84 # 4. Month-by-month monitoring
85 # -----
86 reference_feature = X_train[:, 0] # monitor feature 0
87 monitoring = []
88
89 for m in months:
90     X_m, y_m = monthly_data[m]
91     proba_m = pipeline.predict_proba(X_m)[:, 1]
92
93     # PSI on feature 0
94     psi_f0 = compute_psi(reference_feature, X_m[:, 0])
95
96     # KS test on feature 0 (alternative to PSI)
97     ks_stat, ks_pval = stats.ks_2samp(reference_feature, X_m[:, 0])
98
99     # AUC requires ground truth labels (only available with delay)
100    auc = roc_auc_score(y_m, proba_m)
101
102    # Score distribution PSI (model output drift)
103    ref_scores = pipeline.predict_proba(X_train[:300])[:, 1]
104    psi_score = compute_psi(ref_scores, proba_m)
105
106    monitoring.append({
107        'month': m,
108        'psi_feature_0': psi_f0,
109        'ks_statistic' : ks_stat,
110        'ks_pvalue'    : ks_pval,
111        'psi_score'     : psi_score,
112        'auc'          : auc,
113    })
114
115    mon_df = pd.DataFrame(monitoring).set_index('month')
116
117    print("Monthly Monitoring Dashboard:")
118    print(mon_df.round(4).to_string())
119
120 # -----
121 # 5. Monitoring dashboard plot
122 # -----
123 fig, axes = plt.subplots(2, 2, figsize=(13, 9))
124
125 # PSI feature
126 ax = axes[0,0]
127 colors = ['#e65100' if v >= 0.2 else '#f57c00' if v >= 0.1
128           else '#1b5e20' for v in mon_df['psi_feature_0']]

```

```

129 ax.bar(mon_df.index, mon_df['psi_feature_0'], color=colors, alpha=0.
    85)
130 ax.axhline(0.1, color='orange', linestyle='--', linewidth=1.5, label
    ='PSI=0.1 (caution)')
131 ax.axhline(0.2, color='red',      linestyle='--', linewidth=1.5, label
    ='PSI=0.2 (alert)')
132 ax.set_title("PSI Feature 0 (Data Drift)", fontweight='bold')
133 ax.set_xlabel("Month"); ax.set_ylabel("PSI")
134 ax.legend(fontsize=8); ax.grid(alpha=0.3)
135 ax.set_xticks(months)
136
137 # KS statistic
138 ax = axes[0,1]
139 ax.plot(mon_df.index, mon_df['ks_statistic'],
140         marker='o', color='#0d47a1', linewidth=2)
141 ax.axhline(0.1, color='red', linestyle='--', label='KS=0.1 (alert)')
142 ax.fill_between(mon_df.index, 0, mon_df['ks_statistic'],
143                 alpha=0.15, color='#0d47a1')
144 ax.set_title("KS Statistic Feature 0", fontweight='bold')
145 ax.set_xlabel("Month"); ax.set_ylabel("KS Statistic")
146 ax.legend(fontsize=8); ax.grid(alpha=0.3)
147 ax.set_xticks(months)
148
149 # Score distribution PSI
150 ax = axes[1,0]
151 ax.bar(mon_df.index, mon_df['psi_score'], color='#00838f', alpha=0.
    85)
152 ax.axhline(0.1, color='orange', linestyle='--', linewidth=1.5)
153 ax.axhline(0.2, color='red',      linestyle='--', linewidth=1.5)
154 ax.set_title("PSI Model Score Distribution", fontweight='bold')
155 ax.set_xlabel("Month"); ax.set_ylabel("PSI")
156 ax.grid(alpha=0.3); ax.set_xticks(months)
157
158 # AUC over time
159 ax = axes[1,1]
160 ax.plot(mon_df.index, mon_df['auc'],
161         marker='s', color='#1b5e20', linewidth=2.5, markersize=8)
162 ax.fill_between(mon_df.index, mon_df['auc'].iloc[0] - 0.03,
163                 mon_df['auc'], alpha=0.2, color='#1b5e20')
164 ax.axhline(mon_df['auc'].iloc[0] * 0.95, color='red', linestyle='--',
165            ,
166            label='5% degradation threshold')
167 ax.set_title("Model AUC Over Time", fontweight='bold')
168 ax.set_xlabel("Month"); ax.set_ylabel("AUC")
169 ax.legend(fontsize=9); ax.grid(alpha=0.3)
170 ax.set_xticks(months)
171
172 plt.suptitle("Production Monitoring Dashboard", fontsize=14,
173             fontweight='bold', y=1.01)
174 plt.tight_layout()
175 plt.savefig("monitoring_dashboard.png", dpi=150, bbox_inches='tight'
176            )
177 plt.close()
178 print("\nMonitoring dashboard saved.")

```

The monitoring dashboard separates the *early warning signals* (PSI on raw features, visible immediately in production) from the *lagging performance signal* (AUC, which requires

ground truth labels that may arrive weeks later). The PSI on the model's output score distribution is the most valuable early warning: it detects shifts in how the model is scoring the population even before labels arrive.

In a real deployment, true labels are almost never available immediately. Monitor the model's score distribution using PSI on a weekly or daily cadence. When $\text{PSI} \geq 0.2$, trigger a manual review. When $\text{PSI} \geq 0.2$ and labels arrive confirming AUC degradation, trigger retraining. Build this alerting into your MLOps pipeline from day one.

12 Learning Curves & Diagnosing Underfitting vs Overfitting

12.1 The Bias-Variance Trade-off in Practice

Every model lives on a spectrum between two failure modes.

High bias (underfitting): the model is too simple to capture the patterns in the data. Both training error and test error are high. The model makes systematic errors — the same type of mistake regardless of how much data it sees. Solution: increase model complexity or engineer better features.

High variance (overfitting): the model is too complex and has memorised the training data. Training error is low but test error is high. Performance improves with more data. Solution: regularisation, dropout, pruning, or more data.

Learning curves are the primary tool for diagnosing which failure mode dominates.

12.2 Reading Learning Curves

A **learning curve** plots training and validation performance as a function of training set size. The diagnosis is visual:

- **Both curves converge to a high error plateau:** high bias. More data will not help. The model needs to be more complex.
- **Large gap between training and validation curves:** high variance. The model would benefit from more data. If data cannot be obtained, apply regularisation.
- **Curves converge to a low error:** good fit. The model has learned the signal without overfitting.

Learning Curves and Complexity Curves

```
1 import numpy as np
2 import matplotlib
3 matplotlib.use('Agg')
4 import matplotlib.pyplot as plt
5 from sklearn.datasets import make_classification
6 from sklearn.ensemble import GradientBoostingClassifier
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.model_selection import learning_curve, validation_curve
10 , StratifiedKFold
11 from sklearn.pipeline import Pipeline
12 from sklearn.preprocessing import StandardScaler
13
14 X, y = make_classification(
15     n_samples=3000, n_features=20, n_informative=8,
16     n_redundant=4, weights=[0.65, 0.35], random_state=42
17 )
18 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
```

```

19
20 # -----
21 # 1. Learning curves for three models
22 #   train_sizes: fractions of training data to use in each step.
23 #   Returns arrays of shape (n_steps, n_cv_folds).
24 # -----
25 models = {
26     'Logistic Regression\n(High Bias)': Pipeline([
27         ('sc', StandardScaler()),
28         ('clf', LogisticRegression(C=0.001, max_iter=1000)) # over-
                regularised
29     ]),
30     'Gradient Boosting\n(Good Fit)': Pipeline([
31         ('sc', StandardScaler()),
32         ('clf', GradientBoostingClassifier(
33             n_estimators=100, max_depth=3,
34             learning_rate=0.1, random_state=42
35         ))
36     ]),
37     'Decision Tree\n(High Variance)': Pipeline([
38         ('sc', StandardScaler()),
39         ('clf', DecisionTreeClassifier(max_depth=None, random_state=
40             42))
41     ]),
42 }
43 train_sizes = np.linspace(0.05, 1.0, 15)
44
45 fig, axes = plt.subplots(1, 3, figsize=(15, 5), sharey=True)
46
47 for ax, (name, pipe) in zip(axes, models.items()):
48     sizes, train_scores, val_scores = learning_curve(
49         pipe, X, y,
50         train_sizes=train_sizes,
51         cv=cv,
52         scoring='roc_auc',
53         n_jobs=-1,
54         return_times=False
55     )
56
57     # Mean and std across folds
58     train_mean = train_scores.mean(axis=1)
59     train_std = train_scores.std(axis=1)
60     val_mean = val_scores.mean(axis=1)
61     val_std = val_scores.std(axis=1)
62
63     ax.plot(sizes, train_mean, 'o-', color='#0d47a1', linewidth=2,
64             markersize=5, label='Training AUC')
65     ax.fill_between(sizes, train_mean - train_std,
66                    train_mean + train_std, alpha=0.15, color='#0
67                        d47a1')
68
69     ax.plot(sizes, val_mean, 's-', color='#e65100', linewidth=2,
70             markersize=5, label='Validation AUC')
71     ax.fill_between(sizes, val_mean - val_std,
72                    val_mean + val_std, alpha=0.15, color='#e65100')
73
74     gap = (train_mean[-1] - val_mean[-1])

```



```

74     ax.set_title(f"{name}\nFinal gap: {gap:.3f}", fontsize=11,
75                 fontweight='bold')
76     ax.set_xlabel("Training set size", fontsize=10)
77     ax.legend(fontsize=9)
78     ax.grid(alpha=0.3)
79     ax.set_ylim(0.4, 1.05)
80
81 axes[0].set_ylabel("AUC", fontsize=11)
82 plt.suptitle("Learning Curves: High Bias vs Good Fit vs High
83             Variance",
84             fontsize=13, fontweight='bold', y=1.02)
85 plt.tight_layout()
86 plt.savefig("learning_curves.png", dpi=150, bbox_inches='tight')
87 plt.close()
88
89 # -----
90 # 2. Validation curve (Complexity Curve)
91 # Vary a single hyperparameter and observe how training and
92 # validation performance change. This reveals the sweet spot
93 # between underfitting and overfitting for that parameter.
94 # -----
95
96 param_range = np.arange(1, 10) # max_depth from 1 to 9
97
98 pipe_dt = Pipeline([
99     ('sc', StandardScaler()),
100     ('clf', DecisionTreeClassifier(random_state=42))
101 ])
102
103 train_scores_vc, val_scores_vc = validation_curve(
104     pipe_dt, X, y,
105     param_name='clf__max_depth',
106     param_range=param_range,
107     cv=cv,
108     scoring='roc_auc',
109     n_jobs=-1
110 )
111
112 fig, ax = plt.subplots(figsize=(9, 5))
113
114 tr_m = train_scores_vc.mean(axis=1)
115 tr_s = train_scores_vc.std(axis=1)
116 va_m = val_scores_vc.mean(axis=1)
117 va_s = val_scores_vc.std(axis=1)
118
119 ax.plot(param_range, tr_m, 'o-', color='#0d47a1', lw=2, label='
120         Training AUC')
121 ax.fill_between(param_range, tr_m - tr_s, tr_m + tr_s,
122                alpha=0.15, color='#0d47a1')
123 ax.plot(param_range, va_m, 's-', color='#e65100', lw=2, label='
124         Validation AUC')
125 ax.fill_between(param_range, va_m - va_s, va_m + va_s,
126                alpha=0.15, color='#e65100')
127
128 best_depth = param_range[np.argmax(va_m)]
129 ax.axvline(best_depth, color='green', linestyle=':', lw=1.5,
130            label=f'Best depth={best_depth}')
131
132 ax.set_xlabel("max_depth", fontsize=12)

```

```

129 ax.set_ylabel("AUC", fontsize=12)
130 ax.set_title("Validation Curve Decision Tree max_depth\n"
131             "(Left=underfit, Right=overfit, Peak=sweet spot)",
132             fontsize=12, fontweight='bold')
133 ax.legend(fontsize=10)
134 ax.grid(alpha=0.3)
135 ax.set_xticks(param_range)
136 plt.tight_layout()
137 plt.savefig("validation_curve.png", dpi=150, bbox_inches='tight')
138 plt.close()
139 print("Learning and validation curves saved.")
140
141 # -----
142 # 3. Diagnosis summary
143 # -----
144 print("\n--- Bias-Variance Diagnosis ---")
145 for name, pipe in models.items():
146     sizes, tr, val = learning_curve(
147         pipe, X, y, train_sizes=[0.8, 1.0],
148         cv=cv, scoring='roc_auc', n_jobs=-1
149     )
150     tr_final = tr[-1].mean()
151     val_final = val[-1].mean()
152     gap = tr_final - val_final
153     if tr_final < 0.75 and val_final < 0.75:
154         diagnosis = "HIGH BIAS (underfit)"
155     elif gap > 0.10:
156         diagnosis = "HIGH VARIANCE (overfit)"
157     else:
158         diagnosis = "Good fit"
159     clean = name.replace('\n', ' ')
160     print(f" {clean:<35}: train={tr_final:.3f} val={val_final:.3f}
161           "
162           f"gap={gap:.3f} -> {diagnosis}")

```

12.3 The Validation Curve: Finding the Complexity Sweet Spot

The *validation curve* is the natural complement to the learning curve. Instead of varying the amount of data, it varies a single hyperparameter controlling complexity. The training AUC always increases with complexity; the validation AUC rises, peaks, and then falls as the model begins to overfit. The peak of the validation curve is the optimal complexity setting — this is what hyperparameter tuning is searching for, but the validation curve makes it visually explicit.

13 Multiclass & Multi-label Metrics

13.1 Beyond Binary Classification

Many real problems involve more than two classes: document categorisation, product defect type classification, disease subtype identification, sentiment with multiple levels. The metrics from binary classification extend to multiclass settings, but the averaging strategy becomes a critical decision that changes what is being measured.

13.2 Averaging Strategies

Given per-class precision P_c , recall R_c , and F_1^c for each class $c \in \{1, \dots, K\}$:

Macro averaging: unweighted mean across classes. Every class contributes equally regardless of size. Use this when you care equally about rare and common classes.

$$F_1^{\text{macro}} = \frac{1}{K} \sum_{c=1}^K F_1^c$$

Weighted averaging: weighted mean by class support (number of true instances). Reflects the contribution of each class proportionally. Use this when larger classes matter more.

$$F_1^{\text{weighted}} = \frac{1}{\sum_c n_c} \sum_{c=1}^K n_c \cdot F_1^c$$

Micro averaging: aggregates TP, FP, FN across all classes before computing. For precision and recall, micro-average equals accuracy when all predictions are considered.

13.3 Cohen's Kappa: Agreement Beyond Chance

Cohen's Kappa measures agreement between predicted and actual labels, corrected for what would be expected by chance. It ranges from -1 (perfect disagreement) through 0 (chance agreement) to 1 (perfect agreement):

$$\kappa = \frac{p_o - p_e}{1 - p_e}$$

where p_o is the observed accuracy and p_e is the expected accuracy under chance (computed from marginal distributions). Kappa is especially useful for imbalanced multiclass problems where accuracy is misleading.

Comprehensive Multiclass Evaluation

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib
4 matplotlib.use('Agg')
5 import matplotlib.pyplot as plt
```

```

6  import seaborn as sns
7  from sklearn.datasets import make_classification
8  from sklearn.ensemble import GradientBoostingClassifier
9  from sklearn.model_selection import StratifiedKFold, cross_val_
   predict
10 from sklearn.pipeline import Pipeline
11 from sklearn.preprocessing import StandardScaler, label_binarize
12 from sklearn.metrics import (
13     classification_report, confusion_matrix,
14     cohen_kappa_score, roc_auc_score,
15     f1_score, ConfusionMatrixDisplay
16 )
17
18 # -----
19 # 1. 4-class dataset: imbalanced (mirrors real document
   categorisation)
20 # -----
21 X, y = make_classification(
22     n_samples=3000, n_features=20, n_informative=12,
23     n_classes=4, n_clusters_per_class=1,
24     weights=[0.4, 0.3, 0.2, 0.1], # class 3 is rare
25     random_state=42
26 )
27 class_names = ['Class A', 'Class B', 'Class C', 'Class D']
28 print(f"Class distribution: {np.bincount(y)}")
29
30 # -----
31 # 2. Model and out-of-fold predictions
32 # -----
33 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
34 pipeline = Pipeline([
35     ('scaler', StandardScaler()),
36     ('clf', GradientBoostingClassifier(
37         n_estimators=200, max_depth=4,
38         learning_rate=0.05, random_state=42
39     ))
40 ])
41
42 oof_labels = cross_val_predict(pipeline, X, y, cv=cv, n_jobs=-1)
43 oof_proba = cross_val_predict(pipeline, X, y, cv=cv,
44                               method='predict_proba', n_jobs=-1)
45
46 # -----
47 # 3. Full classification report (per-class precision/recall/F1)
48 # -----
49 print("\nClassification Report:")
50 print(classification_report(y, oof_labels, target_names=class_names)
51 )
52
53 # -----
54 # 4. Aggregate metrics with different averaging strategies
55 # -----
56 for avg in ['macro', 'weighted', 'micro']:
57     f1 = f1_score(y, oof_labels, average=avg)
58     print(f"F1 ({avg:<8}): {f1:.4f}")
59
60 kappa = cohen_kappa_score(y, oof_labels)
61 print(f"Cohen's Kappa : {kappa:.4f} ")

```

```

61         f"({'good' if kappa > 0.6 else 'moderate' if kappa > 0.4 else
           'poor'})")
62
63 # -----
64 # 5. Multiclass ROC-AUC (One-vs-Rest)
65 # For multiclass AUC, we treat each class as binary (one vs. rest
66 # and average. macro averages class AUCs equally; weighted by
   support.
67 # -----
68 auc_macro = roc_auc_score(y, oof_proba, multi_class='ovr',
   average='macro')
69 auc_weighted = roc_auc_score(y, oof_proba, multi_class='ovr',
   average='weighted')
70 print(f"AUC-OVR (macro) : {auc_macro:.4f}")
71 print(f"AUC-OVR (weighted): {auc_weighted:.4f}")
72
73 # -----
74 # 6. Normalised confusion matrix
75 # Normalising by true labels (rows) shows recall per class.
76 # Normalising by predicted labels (cols) shows precision per
   class.
77 # -----
78 fig, axes = plt.subplots(1, 2, figsize=(13, 5))
79
80 cm_raw = confusion_matrix(y, oof_labels)
81 cm_norm = confusion_matrix(y, oof_labels, normalize='true')
82
83 for ax, cm, title, fmt in zip(
84     axes,
85     [cm_raw, cm_norm],
86     ['Confusion Matrix (counts)', 'Confusion Matrix (recall-
       normalised)'],
87     ['d', '.2f']):
88     disp = ConfusionMatrixDisplay(confusion_matrix=cm,
       display_labels=class_names)
89     disp.plot(ax=ax, colorbar=False, cmap='Blues', values_format=fmt
90     )
91     ax.set_title(title, fontsize=12, fontweight='bold')
92     ax.tick_params(axis='x', rotation=20)
93
94 plt.tight_layout()
95 plt.savefig("multiclass_confusion.png", dpi=150, bbox_inches='tight'
96 )
97 plt.close()
98
99 # -----
100 # 7. Per-class metrics visualisation
101 # -----
102 report = classification_report(y, oof_labels,
103                               target_names=class_names, output_
       dict=True)
104 metrics_df = pd.DataFrame({
105     cls: {
106         'Precision': report[cls]['precision'],
107         'Recall' : report[cls]['recall'],
108         'F1' : report[cls]['f1-score'],

```

```

109         'Support' : report[cls]['support'],
110     }
111     for cls in class_names
112 }).T
113
114 fig, ax = plt.subplots(figsize=(9, 5))
115 x = np.arange(len(class_names))
116 w = 0.25
117
118 ax.bar(x - w, metrics_df['Precision'], w, label='Precision',
119        color='#0d47a1', alpha=0.85)
120 ax.bar(x, metrics_df['Recall'], w, label='Recall',
121        color='#00838f', alpha=0.85)
122 ax.bar(x + w, metrics_df['F1'], w, label='F1 Score',
123        color='#e65100', alpha=0.85)
124
125 ax.set_xticks(x)
126 ax.set_xticklabels(
127     [f"{c}\n(n={int(metrics_df.loc[c, 'Support'])})" for c in class_
128      names],
129     fontsize=10
130 )
131 ax.set_ylim(0, 1.1)
132 ax.set_ylabel("Score", fontsize=11)
133 ax.set_title(f"Per-class Metrics (Kappa={kappa:.3f})",
134             fontsize=12, fontweight='bold')
135 ax.legend(fontsize=10)
136 ax.grid(axis='y', alpha=0.3)
137 plt.tight_layout()
138 plt.savefig("multiclass_perclass.png", dpi=150, bbox_inches='tight')
139 plt.close()
140 print("\nMulticlass visualisations saved.")
141
142 # -----
143 # 8. Which averaging strategy to report?
144 # -----
145 print("\nAveraging Strategy Guide:")
146 print(f"    macro    F1={f1_score(y, oof_labels, average='macro'):.4f}"
147       "
148       "-- equal weight to all classes (good for rare class
149       monitoring)")
150 print(f"    weighted F1={f1_score(y, oof_labels, average='weighted'):.4f}"
151       "
152       "-- weight by support (good for overall performance)")
153 print(f"    micro    F1={f1_score(y, oof_labels, average='micro'):.4f}"
154       "
155       "-- global TP/FP/FN (equals accuracy when all samples counted)"
156       "
157       Kappa={kappa:.4f} -- corrected for chance (robust to
158       imbalance)")

```

13.4 The Normalised Confusion Matrix as a Debugging Tool

The raw confusion matrix shows counts; the normalised version (row-normalised by true class) shows recall per class. For a 4-class problem, reading down each column of the normalised matrix reveals which classes are being confused with each other. If Class

D (rare) is being predicted as Class A 60% of the time, the model has not learned a discriminative representation of Class D — a signal to collect more examples, oversample, or apply class-weighted loss.

Situation	Recommended Metric	Averaging
All classes equally important, balanced	Accuracy or Macro F1	Macro
Rare class matters most	Macro F1, per-class recall	Macro
Overall performance for a mixed audience	Weighted F1	Weighted
Comparing to human annotators	Cohen's Kappa	N/A
Probabilistic ranking across classes	Multiclass (OVR)	AUC Macro/Weighted

Table 6: Metric selection guide for multiclass problems.

Regression Models

- ☐ Reported MAE *and* RMSE (sensitivity to outliers differs).
- ☐ R^2 reported alongside residual plots — not in isolation.
- ☐ Four-panel residual diagnostic completed (patterns indicate missing structure).
- ☐ MAPE or sMAPE used only when target values are not near zero.

Calibration

- ☐ If probabilities are used directly (not just for ranking): reliability diagram inspected.
- ☐ Brier score reported alongside AUC.
- ☐ Calibration applied (Platt or isotonic) if reliability curve deviates from diagonal.

Statistical Testing

- ☐ Model comparisons confirmed with corrected resampled t-test or McNemar's test.
- ☐ Bonferroni correction applied when comparing more than two models.
- ☐ Difference in CV means reported with confidence interval, not just point estimate.

Fairness & Bias

- ☐ Protected attributes identified and excluded from model training.

- ☐ Per-group metrics computed for all relevant protected groups.
- ☐ Disparate Impact Ratio checked (≥ 0.8 required in most jurisdictions).
- ☐ Chosen fairness criterion (demographic parity vs equalized odds) documented and justified.

Production Monitoring

- ☐ PSI monitoring pipeline built before deployment, not after.
- ☐ Score distribution PSI and feature PSI monitored separately.
- ☐ Retraining triggers defined (e.g., $\text{PSI} \geq 0.2$ or $\text{AUC degradation} > 5\%$).
- ☐ Baseline reference distributions saved at deployment time.

Bias-Variance and Learning Curves

- ☐ Learning curves generated and inspected before finalising model complexity.
- ☐ Validation curve run for at least one key hyperparameter.
- ☐ Diagnosis (high bias / high variance / good fit) documented and acted upon.

Multiclass

- ☐ Averaging strategy (macro/weighted/micro) chosen based on class balance and business importance.
- ☐ Per-class precision, recall, and F1 inspected — not just aggregate metrics.
- ☐ Normalised confusion matrix inspected for inter-class confusions.
- ☐ Cohen's Kappa reported when comparing to human baseline or annotator agreement.

Appendix: Quick Reference

Key Formulas

$$\text{Precision} = \frac{TP}{TP + FP} \quad (1)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2)$$

$$F_\beta = (1 + \beta^2) \cdot \frac{P \cdot R}{\beta^2 P + R} \quad (3)$$

$$\text{Balanced Accuracy} = \frac{1}{2} \left(\frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right) \quad (4)$$

$$\text{SHAP: } f(\mathbf{x}) = \mathbb{E}[f(X)] + \sum_j \phi_j \quad (5)$$

$$\text{PDP: } \hat{f}_j(x_j) = \frac{1}{n} \sum_{i=1}^n \hat{f}(x_j, \mathbf{x}_{-j}^{(i)}) \quad (6)$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (7)$$

$$\text{Brier Score} = \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i)^2 \quad (8)$$

$$\text{PSI} = \sum_b (q_b - p_b) \ln(q_b/p_b) \quad (9)$$

$$\kappa = \frac{p_o - p_e}{1 - p_e} \quad (10)$$

Library Quick-Install

Required packages

```
pip install scikit-learn shap lime scikit-optimize scipy matplotlib
pandas numpy seaborn
```

Key Classes and Functions

Class / Function	Purpose
StratifiedKFold	K-fold with class distribution preserved
TimeSeriesSplit	Walk-forward CV for temporal data
cross_validate	Multi-metric CV with train/test scores
cross_val_predict	Out-of-fold predictions for full-dataset evaluation
Pipeline	Chain preprocessing + model (prevents leakage)
GridSearchCV	Exhaustive hyperparameter search
RandomizedSearchCV	Stochastic hyperparameter search
BayesSearchCV	Bayesian hyperparameter optimisation (skopt)
average_precision_score	Area under precision-recall curve
balanced_accuracy_score	Accuracy corrected for class imbalance
precision_recall_curve	Threshold-free P-R analysis
shap.TreeExplainer	Fast SHAP for tree models
shap.summary_plot	Global SHAP visualisation
LimeTabularExplainer	Local LIME explanation for tabular data
PartialDependenceDisplay	PDP and ICE plots
mean_absolute_error	MAE for regression
r2_score	Coefficient of determination (R^2)
CalibratedClassifierCV	Platt scaling / isotonic regression
calibration_curve	Reliability diagram data
brier_score_loss	Probability calibration quality
learning_curve	Train/val score vs training set size
validation_curve	Train/val score vs hyperparameter value
cohen_kappa_score	Agreement metric for multiclass
roc_auc_score	AUC-ROC, supports multiclass OVR
ConfusionMatrixDisplay	Normalised multiclass confusion matrix

Acknowledgments



Esmail Rezaei, Ph.D.

Thank you for reading this guide! I hope this tutorial has helped you understand how to preprocess text data for transformer models. If you found this useful, I would greatly appreciate your support by connecting with me on my professional networks.

GitHub
esmaeil-rezaei

LinkedIn
esmaeil-rezaei

Feel free to explore my repositories on GitHub for more machine learning projects and connect with me on LinkedIn to stay updated on my latest work. Your feedback and suggestions are always welcome!