

Professional Certificate in Risk Modeling with Machine Learning

## Model Evaluation And Selection

**AIC (Akaike Information Criterion)** – related terms: BIC, model complexity, likelihood. AIC estimates the relative quality of statistical models for a given dataset. It balances model fit (via the log-likelihood) against the number of estimated parameters, penalising excessive complexity. Lower AIC values indicate a preferable model. Example: Two logistic-regression models for credit-default prediction have AIC scores of 1245 and 1260; the model with 1245 is selected. Practical application: In risk modeling, AIC helps compare alternative specifications (e.G., Adding interaction terms) without over-fitting. Challenges: AIC is an asymptotic approximation; with small samples it may favour overly complex models. It also does not provide an absolute measure of goodness-of-fit, only a relative ranking.

**AUC-ROC (Area Under the Receiver Operating Characteristic Curve)** – related terms: ROC curve, discrimination, c-statistic. AUC-ROC quantifies a classifier’s ability to rank positive instances higher than negative ones across all possible decision thresholds. An AUC of 0.5 Denotes random guessing; 1.0 Denotes perfect discrimination. Example: A fraud-detection model yields an AUC of 0.87, Meaning that in 87 % of randomly chosen fraud-non-fraud pairs, the model assigns a higher probability to the fraud case. Practical application: Regulators often require high discrimination for credit-risk scores; AUC-ROC is a standard reporting metric. Challenges: AUC is insensitive to calibration; two models with identical AUC can have very different probability estimates. It also masks performance on specific operating points that may be business-critical.

**Bias-Variance Trade-off** – related terms: Overfitting, underfitting, generalisation error. The bias-variance trade-off describes how model error decomposes into bias (error from erroneous assumptions) and variance (error from sensitivity to training data fluctuations). High-bias models are overly simple, while high-variance models capture noise. Example: A shallow decision tree may underfit (high bias) a loan-loss dataset, whereas a deep tree may overfit (high variance). Practical application: Selecting the appropriate depth or regularisation strength for gradient-boosted trees in insurance-pricing models requires balancing bias and variance. Challenges: Quantifying bias and variance directly is difficult; practitioners rely on cross-validation results, which can be noisy for small datasets.

**Calibration Curve** – related terms: Reliability diagram, probability calibration, Brier score. A calibration curve plots predicted probabilities against observed event frequencies, typically in decile bins. Perfect calibration lies on the 45-degree line. Example: A credit-risk model predicts a 10 % default probability for a group of borrowers; the observed default rate is 18 %, indicating under-prediction. Practical application: In insurance, calibrated probability estimates are essential for setting premiums that reflect true risk levels. Challenges: Calibration may deteriorate when models are retrained on new data distributions; recalibrating (e.G., Via Platt scaling) adds complexity and may reduce discrimination.

**Confusion Matrix** – related terms: True positive, false negative, precision, recall. A confusion matrix tabulates the counts of true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN) for a binary classifier at a chosen threshold. It provides the foundation for derived metrics such as accuracy,

precision, recall, and F1 score. Example: For a fraud-detection system, the matrix shows TP = 120, FP = 30, FN = 15, TN = 2835. Practical application: In regulatory reporting, the matrix helps illustrate the trade-off between catching fraud (TP) and inconveniencing legitimate customers (FP). Challenges: The matrix is threshold-dependent; selecting a threshold that satisfies business constraints (e.G., Cost of false alarms) often requires additional analysis.

Cost-Sensitive Learning – related terms: Misclassification cost matrix, utility, threshold optimisation. Cost-sensitive learning incorporates asymmetric costs of errors directly into model training or decision making. Instead of treating all mistakes equally, it weights false positives and false negatives according to business impact. Example: In credit-risk scoring, a false negative (approving a high-risk borrower) may cost \$10 000, while a false positive (rejecting a low-risk applicant) costs \$500 in lost revenue. Practical application: By embedding these costs, a logistic-regression model can be trained to minimise expected loss rather than simple error rate, leading to more profitable decision thresholds. Challenges: Accurate cost quantification is difficult; costs may vary over time, and overly aggressive cost weighting can cause model instability.

Cross-Validation – related terms: k-fold, holdout set, model validation. Cross-validation partitions the data into multiple training and validation folds to assess model performance more robustly than a single holdout split. Common variants include k-fold, stratified k-fold, and leave-one-out. Example: A 5-fold cross-validation of a random-forest credit-risk model yields average AUC = 0.84 With a standard deviation of 0.02. Practical application: In the development of an operational risk model, cross-validation helps detect overfitting before deployment. Challenges: Computational cost grows with the number of folds, especially for complex models. Data leakage (e.G., Preprocessing on the full dataset) can inflate performance estimates.

Decision Threshold – related terms: ROC curve, cost-sensitive learning, operating point. A decision threshold converts predicted probabilities into binary class labels. Moving the threshold shifts the balance between TP and FP rates. Example: Lowering the threshold from 0.5 To 0.3 In a fraud detector increases TP from 120 to 150 but also raises FP from 30 to 70. Practical application: Business units often set thresholds based on acceptable false-positive rates (e.G., Maximum 2% of legitimate transactions flagged). Challenges: The optimal threshold depends on context-specific costs and may change as the underlying data distribution evolves.

F1 Score – related terms: Precision, recall, harmonic mean. The F1 score combines precision ( $TP/(TP+FP)$ ) and recall ( $TP/(TP+FN)$ ) into a single metric via their harmonic mean:  $2 \cdot (\text{Precision} \cdot \text{recall}) / (\text{precision} + \text{recall})$ . It is especially useful when the class distribution is imbalanced. Example: A model with precision = 0.80 And recall = 0.60 Yields an F1 of 0.69. Practical application: In anti-money-laundering screening, where fraudulent cases are rare, the F1 score provides a balanced view of detection capability. Challenges: F1 ignores true negatives, which may be important in certain risk contexts; it also treats precision and recall equally, which may not reflect business priorities.

Gini Coefficient – related terms: AUC-ROC, Lorenz curve, discrimination. The Gini coefficient is a scaled version of AUC, calculated as  $2 \cdot \text{AUC} - 1$ . It ranges from 0 (no discrimination) to 1 (perfect discrimination). Example: An AUC of 0.78 Translates to a Gini of 0.56, Indicating moderate predictive power for a

mortgage-default model. Practical application: Many financial institutions report Gini to regulators as a concise measure of model discrimination. Challenges: Like AUC, Gini does not assess calibration; two models with identical Gini can have divergent probability estimates.

Holdout Validation – related terms: Train-test split, validation set, temporal split. Holdout validation reserves a portion of the data (commonly 20-30%) as a test set that remains untouched during model training and hyper-parameter tuning. Performance on this set provides an unbiased estimate of out-of-sample accuracy. Example: A risk-score model is trained on data up to 2022 and evaluated on 2023 transactions held out as the test set. Practical application: In production pipelines, a holdout set can be refreshed periodically to monitor model drift. Challenges: With limited data, a single holdout split can produce high variance in performance estimates; random splits may also break temporal dependencies important in risk modeling.

K-Fold Cross-Validation – related terms: Cross-validation, stratified sampling, model stability. In k-fold cross-validation, the dataset is divided into k equally sized folds; each fold serves as a validation set once while the remaining k – 1 folds train the model. The average performance across folds is reported. Example: A 10-fold cross-validation of a gradient-boosted model yields mean log-loss = 0.215 With a 95% confidence interval of [0.208, 0.222]. Practical application: Hyper-parameter optimisation frameworks (e.g., Grid search) often use k-fold CV to evaluate each candidate configuration. Challenges: When data exhibit strong temporal autocorrelation, random folding can leak future information into the training folds, overstating performance.

Lift Chart – related terms: Cumulative gains, response rate, targeting efficiency. A lift chart displays the ratio of the model's captured positive events to the baseline capture rate, typically plotted across deciles of predicted risk. Higher lift in early deciles indicates better targeting. Example: In a cross-sell campaign, the top 10% of customers identified by the model generate a lift of 4.5, Meaning they are 4.5 Times more likely to respond than the average customer. Practical application: Insurance underwriters use lift to allocate underwriting resources to high-risk segments efficiently. Challenges: Lift is sensitive to class imbalance and may be misleading if the underlying event rate changes over time.

Log-Loss (Cross-Entropy Loss) – related terms: Likelihood, calibration, negative log-likelihood. Log-loss measures the accuracy of probability predictions; lower values indicate better calibrated and more confident forecasts. For binary outcomes, it is defined as  $-[y \cdot \log(p) + (1-y) \cdot \log(1-p)]$ , averaged over all observations. Example: A model with average log-loss = 0.31 Outperforms a baseline model with log-loss = 0.45 On a credit-default dataset. Practical application: Many Kaggle competitions and regulatory model-validation checklists require reporting log-loss as a primary metric. Challenges: Log-loss heavily penalises overconfident incorrect predictions, which can be problematic when probabilities are estimated from limited data.

Mean Absolute Error (MAE) – related terms: Regression metrics, robustness, L1 loss. MAE computes the average absolute difference between predicted and actual continuous outcomes. It is less sensitive to outliers than Mean Squared Error. Example: A loss-given-default (LGD) model predicts LGD values with MAE = 0.07, Indicating an average absolute deviation of 7 percentage points. Practical application: In loss-provision calculations, MAE provides an intuitive measure of forecast error for regulatory reporting. Challenges: MAE does not penalise larger errors as heavily as MSE, potentially obscuring the impact of

extreme mis-predictions on capital requirements.

Mean Squared Error (MSE) – related terms: RMSE, variance, L2 loss. MSE averages the squared differences between predicted and actual values, giving greater weight to larger errors. It is the natural loss function for ordinary least-squares regression. Example: An asset-price-prediction model achieves MSE = 0.018, implying a root-mean-square error (RMSE) of about 13.4% of the price range. Practical application: In actuarial reserving, MSE helps compare competing severity-distribution fits. Challenges: MSE is highly sensitive to outliers; a few extreme residuals can dominate the metric, leading to misleading model rankings.

Model Overfitting – related terms: High variance, regularisation, training error vs. Validation error. Overfitting occurs when a model captures noise or idiosyncrasies of the training data, resulting in poor generalisation to unseen data. Indicators include a large gap between training performance (e.g., High accuracy) and validation performance. Example: A neural network with 500k parameters achieves 99% training accuracy on a credit-risk dataset but only 71% validation accuracy, signalling overfitting. Practical application: Early-stopping, dropout, and pruning are techniques used to curb overfitting in deep-learning risk models. Challenges: Detecting overfitting early requires reliable validation schemes; in highly imbalanced risk datasets, validation metrics may be noisy, obscuring the symptom.

Model Underfitting – related terms: High bias, insufficient complexity, learning curve. Underfitting arises when a model is too simple to capture the underlying data structure, leading to high errors on both training and validation sets. Example: A linear regression applied to a non-linear loss-severity relationship yields  $R^2 = 0.22$  on both train and test sets, indicating underfitting. Practical application: Adding polynomial features or switching to a more flexible algorithm (e.g., Gradient boosting) can alleviate underfitting in risk-score development. Challenges: Over-compensating for underfitting may swing the model into overfitting territory; balancing the two requires systematic hyper-parameter tuning.

Precision – related terms: Positive predictive value, false discovery rate, confusion matrix. Precision measures the proportion of predicted positives that are true positives:  $TP/(TP+FP)$ . It reflects the model's ability to avoid false alarms. Example: In a fraud detection system, precision = 0.80 means that 80% of flagged transactions are indeed fraudulent. Practical application: High precision is essential when false positives are costly, such as when manual investigation resources are limited. Challenges: Precision alone can be misleading in highly imbalanced settings; a model that predicts "no fraud" for every case achieves 100% precision but zero recall.

Recall (Sensitivity) – related terms: True positive rate, detection rate, confusion matrix. Recall quantifies the proportion of actual positives correctly identified:  $TP/(TP+FN)$ . It captures the model's ability to detect events of interest. Example: A credit-risk model with recall = 0.92 correctly flags 92% of defaulting borrowers. Practical application: In regulatory stress testing, high recall ensures that most high-risk exposures are captured for capital adequacy analysis. Challenges: Maximising recall often raises false-positive rates, increasing operational costs; a balanced approach with precision or cost-sensitive metrics is usually required.

ROC Curve (Receiver Operating Characteristic Curve) – related terms: AUC, trade-off, threshold analysis. The ROC curve plots the true-positive rate (TPR) against the false-positive rate (FPR) for varying decision

thresholds. It visualises the trade-off between sensitivity and specificity across the entire range of possible thresholds. Example: A model's ROC curve bows toward the upper-left corner, indicating strong discrimination, whereas a diagonal line denotes random guessing. Practical application: Risk analysts compare ROC curves of competing models to select the one offering the best overall trade-off before fixing a business-specific threshold. Challenges: ROC curves can be overly optimistic in highly imbalanced datasets; the precision-recall curve may be more informative when the positive class is rare.

Sensitivity Analysis – related terms: Scenario testing, parameter perturbation, model robustness. Sensitivity analysis examines how variations in model inputs or assumptions affect outputs. It helps identify variables that most influence risk estimates. Example: Varying the default probability by  $\pm 10\%$  changes the estimated capital requirement by  $\pm 4\%$ , highlighting sensitivity to PD assumptions. Practical application: In credit-portfolio modelling, sensitivity analysis informs stress-testing procedures required by supervisory bodies. Challenges: High-dimensional models can make exhaustive sensitivity testing computationally prohibitive; surrogate models or sampling techniques (e.g., Sobol indices) are often employed.

Specificity (True Negative Rate) – related terms: False positive rate,  $TN/(TN+FP)$ , confusion matrix. Specificity measures the proportion of actual negatives correctly identified:  $TN/(TN+FP)$ . It reflects the model's ability to correctly reject non-events. Example: An underwriting model with specificity = 0.95 Correctly classifies 95% of low-risk applicants as non-defaults. Practical application: High specificity reduces the number of unnecessary rejections, preserving customer goodwill in retail banking. Challenges: Raising specificity typically lowers sensitivity; the optimal balance depends on the relative costs of false positives versus false negatives.

Training Set – related terms: Fit, learning phase, data split. The training set comprises observations used to estimate model parameters. All preprocessing steps that affect the model (e.g., Scaling, encoding) should be derived from this set to avoid leakage. Example: A random-forest classifier for operational risk is trained on 70% of the historical loss events. Practical application: Maintaining a clean separation between training and validation data is a cornerstone of model governance in financial institutions. Challenges: In time-series risk data, the training set must respect chronological order; random shuffling can introduce look-ahead bias.

Validation Set – related terms: Holdout, model selection, hyper-parameter tuning. The validation set is a subset of data not used for fitting model parameters but for assessing performance during development. It guides hyper-parameter choices and early-stopping decisions. Example: During grid search, a validation set of 15% of the data determines the optimal number of trees for a gradient-boosting model. Practical application: Regulatory model-validation frameworks often require a separate validation set to demonstrate that model performance is not an artifact of over-fitting to the training data. Challenges: When data are scarce, allocating a sizable validation set reduces the effective training sample, potentially degrading model quality; cross-validation can mitigate this but adds computational overhead.

Variance Inflation Factor (VIF) – related terms: Multicollinearity, regression diagnostics, feature selection. VIF quantifies how much the variance of an estimated regression coefficient is inflated due to linear dependence with other predictors.  $VIF = 1$  indicates no correlation; values above 5–10 suggest problematic multicollinearity. Example: In a logistic-regression credit-risk model, the VIF for "total debt" is 12, prompting

the analyst to drop or combine correlated variables. Practical application: Reducing multicollinearity improves interpretability and stabilises coefficient estimates, which is critical for models that must be explained to regulators. Challenges: Removing variables solely based on VIF may discard useful predictive information; domain expertise is required to balance interpretability and predictive power.

Weighted Accuracy – related terms: Class weighting, cost matrix, imbalanced classification. Weighted accuracy assigns different importance to correctly classifying each class, often by multiplying each correct prediction by a class-specific weight before averaging. Example: In a fraud detection task, true positives are weighted 5 times more than true negatives, yielding a weighted accuracy of 0.87 Despite an overall accuracy of 0.92. Practical application: Weighted accuracy aligns model evaluation with business objectives when the minority class carries disproportionate risk. Challenges: Choosing appropriate weights requires quantifying the relative economic impact of each error type; mis-specified weights can bias model selection toward undesirable behaviours.

Yield Curve Adjustment – related terms: Macro-stress testing, scenario analysis, risk factor modeling. Yield curve adjustment refers to the transformation of interest-rate risk factors in stress testing to reflect hypothetical shifts (parallel, steepening, flattening) in the term structure. Example: A stress scenario applies a 200-basis-point parallel upward shift to the sovereign yield curve, impacting the valuation of fixed-income portfolios. Practical application: Banks use yield-curve adjustments to estimate potential losses under adverse economic conditions, satisfying supervisory stress-test requirements. Challenges: Selecting realistic shock magnitudes and shapes is non-trivial; overly severe shocks may overstate capital needs, while insufficient shocks may underestimate risk.

Z-Score (Standard Score) – related terms: Standardisation, outlier detection, normalisation. A Z-score measures how many standard deviations an observation lies from the mean of its distribution:  $(X - \mu)/\sigma$ . It is commonly used to standardise features before modelling. Example: A borrower's debt-to-income ratio of 0.45 Yields a Z-score of +1.2 Relative to the industry mean, indicating higher leverage. Practical application: Standardised variables enable algorithms that assume centred data (e.G., Regularised logistic regression) to converge faster and avoid numerical instability. Challenges: Z-scores assume approximate normality; skewed variables may require transformation (e.G., Log) before standardisation to avoid misleading outlier identification.