

Equation Appendix

Notation

- TP, FP, TN, FN – true/false positive/negative counts.
- $N = TP + FP + TN + FN$ – total number of instances.
- $y_i \in \{0, 1\}$ – ground-truth label for instance i .
- $\hat{p}_i \in [0, 1]$ – model-predicted probability of class 1.

1 Binary Classification Metrics

1.1 Accuracy, Precision, Recall

$$\text{Accuracy (ACC)} = \frac{TP + TN}{N}, \quad (1)$$

$$\text{Precision (P)} = \frac{TP}{TP + FP}, \quad (2)$$

$$\text{Recall (R)} = \frac{TP}{TP + FN}. \quad (3)$$

1.2 F₁-score

$$F_1 = 2 \frac{P \cdot R}{P + R}. \quad (4)$$

1.3 Matthews Correlation Coefficient (MCC)

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}. \quad (5)$$

1.4 Area Under the ROC Curve (AUC)

The threshold-independent metric

$$AUC = \int_0^1 TPR(FPR) d(FPR), \quad (6)$$

where $TPR = TP/(TP + FN)$ and $FPR = FP/(FP + TN)$.

2 Loss Functions

2.1 Binary Cross-Entropy (Neural Network)

$$\mathcal{L}_{\text{BCE}} = -\frac{1}{N} \sum_{i=1}^N [y_i \log \hat{p}_i + (1 - y_i) \log(1 - \hat{p}_i)]. \quad (7)$$

2.2 CatBoost Ordered Boosting Objective

CatBoost minimises the additive loss

$$\min_{f \in \mathcal{F}} \sum_{i=1}^N \ell(y_i, f(\mathbf{x}_i)) + \lambda \|f\|_{\text{model}}^2, \quad (8)$$

where ℓ is the log-loss, $\|\cdot\|_{\text{model}}$ the model-specific regulariser and $\lambda > 0$ the L2 term tuned over $[1, 10]$.

3 Log₂ MIC Transformation

$$\text{log2_mic} = \log_2(\text{MIC} / \text{mg L}^{-1}). \quad (9)$$

4 Bootstrapped Confidence Intervals

For an estimator $\hat{\theta}$ and B bootstrap resamples $\{\hat{\theta}^{*(b)}\}_{b=1}^B$:

$$\text{CI}_{95\%} = [\hat{\theta}^{*(0.025)}, \hat{\theta}^{*(0.975)}]. \quad (10)$$

5 DeLong Test for Paired ROC Curves

Given two vectors of AUC contributions V_1 and V_2 :

$$\Delta = \text{AUC}_1 - \text{AUC}_2, \quad (11)$$

$$\sigma_{\Delta}^2 = \frac{\text{Var}(V_1 - V_2)}{N}, \quad (12)$$

$$z = \frac{\Delta}{\sigma_{\Delta}}, \quad p = 2(1 - \Phi(|z|)), \quad (13)$$

where $\Phi(\cdot)$ is the standard normal CDF.

6 Feed-Forward Neural Network Forward Pass

$$\mathbf{h}_1 = \text{ReLU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1), \quad (14)$$

$$\mathbf{h}_2 = \text{ReLU}(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2), \quad (15)$$

$$\hat{y} = \sigma(\mathbf{w}_3^\top \mathbf{h}_2 + b_3), \quad (16)$$

with $\sigma(z) = 1/(1 + e^{-z})$.

7 Bayesian Optimisation Search Space

$$\text{depth} \sim \mathcal{U}\{4, 10\}, \quad (17)$$

$$\eta (\text{learning rate}) \sim \mathcal{U}(0.005, 0.3), \quad (18)$$

$$\lambda_{\text{L2}} \sim \mathcal{U}(1, 10). \quad (19)$$

8 Stratified k -Fold Cross-Validation

Let $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ and strata S_1, \dots, S_K with class proportions preserved; then

$$\text{CV}(\theta) = \frac{1}{K} \sum_{k=1}^K \mathcal{M}(\theta; \mathcal{D}_{\text{train}}^{(k)}, \mathcal{D}_{\text{test}}^{(k)}), \quad (20)$$

where \mathcal{M} is a chosen metric (e.g. AUC).

End of Appendix